

The Role of Data for One-Shot Semantic Segmentation

Timo Lüddecke^{1,§} and Alexander Ecker^{1,2}

¹Institute of Computer Science and Campus Institute Data Science, University of Göttingen

²Max Planck Institute for Dynamics and Self-Organization

§timo.lueddecke@uni-goettingen.de

Due to incorrect data pre-processing, the training set of the LVIS-OneShot dataset contained PASCAL-5ⁱ classes in the original version of this paper, which was published in CVPRW 2021. We updated all affected figures and tables in this correction. For details we refer to the appendix 6.

Abstract

In this work we investigate the potential of larger datasets for one-shot semantic segmentation. While computer vision models are often trained on millions of diverse samples, current one-shot semantic segmentation datasets encompass only a small number of samples (PASCAL-5ⁱ), a small number of classes (PASCAL-5ⁱ and COCO-20ⁱ) or have little variability (FSS-1000). To improve this situation, we introduce LVIS-OneShot, a one-shot variant of the LVIS dataset. With 718 classes and 114,347 images, it exceeds previous datasets substantially in terms of size. By controlled experiments we show that not only the number of images but also the number of different classes is crucial. We analyze transfer learning across common datasets and find that by training on LVIS-OneShot we outperform current state-of-the-art models on PASCAL-5ⁱ. In particular, we observe that a simple baseline model (MaRF) learns to perform one-shot segmentation when trained on a large dataset although it has a generic architecture without strong inductive biases. Code and dataset are available here:

eckerlab.org/code/one-shot-segmentation

1 Introduction

In many fields of computer vision, we witnessed a trend towards training conceptually simple models with a large number of parameters on large-scale datasets involving millions of samples, starting from AlexNet [1] on the ImageNet dataset [2] to the recent Big Transfer [3] and CLIP [4] on text segments from the Internet. Instead of designing mechanisms to solve specific tasks (for example fine-grained recognition), features are learned by scaling up the dataset and applied on downstream tasks by fine-tuning, learning a linear classifier from frozen features or zero-shot transfer [4]. Very recently, similar findings were made for one-shot image classification [5] and object detection [6].

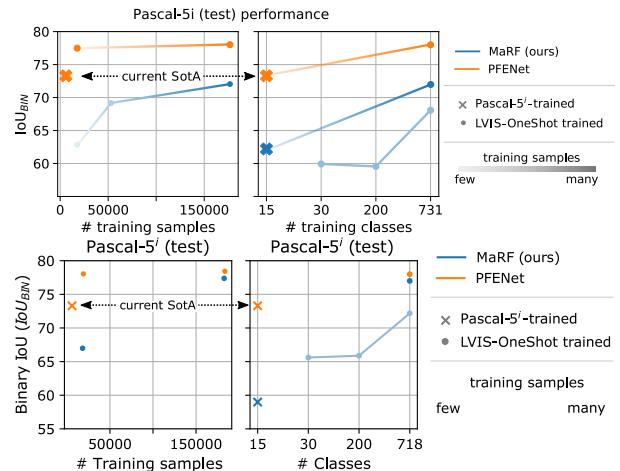


Figure 1: One-shot segmentation on PASCAL-5ⁱ benefits from the larger and more diverse LVIS-OneShot dataset. The simple MaRF baseline profits from more (left) data and more diverse (right) data. Both, PFENet and MaRF exceed the current state-of-the-art (SotA). top new, bottom old (red and orange lines are missing for some reason)

One may wonder to what degree these findings can be replicated for one-shot segmentation, which requires dynamic adaptation at inference time to make pixel-wise predictions. One-shot segmentation models are currently trained and evaluated on datasets with relatively few classes (PASCAL-5ⁱ [7]: 20, COCO-20ⁱ [8]: 80) or little variation in object size and position (FSS-1000 [9]). At the same time, models use customized mechanisms to tackle one-shot segmentation, such as operating on multiple scales [10] or modeling object parts [11, 12, 13] (see Table 1 for an overview). Thus, it is quite possible that scaling up datasets in terms of labels and classes could similarly boost performance in one-shot segmentation and dwarf the advances made by previous attempts at engineering customized architectures for small datasets. In particular, we believe the diversity of labels to

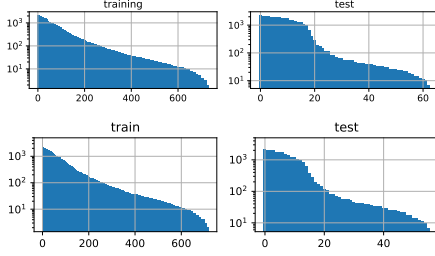


Figure 2: Distribution of the class frequencies in LVIS-OneShot top new, bottom old

play a pivotal role.

In order to investigate these questions, we introduce the LVIS-OneShot dataset which is a variant of the LVIS dataset [14], originally intended for long-tail instance segmentation. Contrary to previous one-shot segmentation datasets, it combines a large number of samples with a large number of classes. Furthermore, we intentionally place all PASCAL-5ⁱ classes into the test set to simplify evaluating performance on PASCAL-5ⁱ. In addition to the dataset, we present the primitive baseline model MaRF without bells and whistles which relies on very simple design: ResNet backbones [15], masked (global) pooling and feature-wise linear modulation (FiLM) [16]. Based on both, LVIS-OneShot dataset and MaRF baseline, as well as available datasets we conduct extensive experiments seeking to understand the role of data in one-shot segmentation.

2 A Diverse Dataset: LVIS-OneShot

With COCO-20ⁱ [20] there already exists a large-scale one-shot (and few-shot) segmentation dataset. However, the small number of classes limits the utility of this dataset and, due to an overlap between PASCAL-5ⁱ and COCO-20ⁱ splits, transfer learning is cumbersome since a new model needs to be trained for each split. We address both shortcomings with a novel dataset (Fig. 2) based on the recently proposed LVIS [14] labels for COCO images [21].

PASCAL-5ⁱ Overlap As a first step, we identify all classes in LVIS that overlap with Pascal. To automate this process, we make use of the WordNet [22] synset assignments of LVIS. First, we manually assign the corresponding

Mechanism	Used by
modeling of object parts	[11, 12]
graph attention on regions/ parts	[13, 17]
multi-scale features / feature pyramid	[10, 18, 17, 13]
high-level feature similarity prior	[10]
metric learning & prototype alignment loss	[19]

Table 1: Overview of mechanisms used for one-shot segmentation.

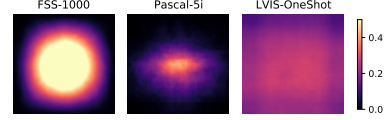


Figure 3: Average images of 5,000 random ground truth segmentations of FSS-1000, PASCAL-5ⁱ and LVIS-OneShot reveal dataset bias.

synset for all 20 Pascal classes. For each LVIS class we recursively traverse the set of hypernyms (i.e. more general meanings) and check if it has an intersection with the Pascal synsets.

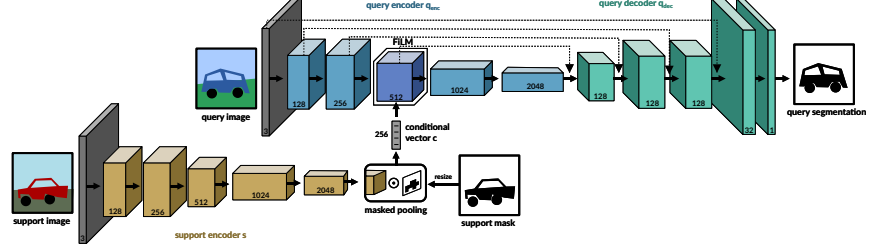
3 Complex and Simple Models

Arguably the best performing model for one-shot segmentation is PFENet [10]. It contains several custom mechanisms, e.g. multi-scale processing and a prior based on high-level feature similarity. Contrary, our baseline model (Fig. 4) relies on only three basic and well-known components: Masked pooling [23], the ResNet architecture [15] and the FiLM conditioning mechanism [16], hence we call it MaRF. It is meant as a low-modelling, generic counterpoint to the complex architecture defined by PFENet. Query and support branch only interact through a single vector.

Support encoding The support encoder s takes the support tuple \mathbf{t} , containing an image with corresponding segmentation, and encodes it into a conditional vector \mathbf{c} . For this, features from a ResNet50 or ResNet18 [15] are extracted at the 3rd or 4th residual block. Then, masked pooling transforms these feature maps into a single vector by averaging all feature vectors that pertain to the object (indicated by the support segmentation map). Given a support image, segmentation pair $\mathbf{t} = (\mathbf{t}_{\text{img}}, \mathbf{t}_{\text{seg}})$ and an operator which resizes \mathbf{t}_{seg} to match the spatial size of the output of s_{enc} , we obtain the conditional vector $\mathbf{c} = \text{AvgPool}(s_{\text{enc}}(\mathbf{t}_{\text{img}}) \odot \text{resize}(\mathbf{t}_{\text{seg}}))$, where \odot is the point-wise multiplication.

Query processing The query network q takes the query image \mathbf{x} and the generated conditional vector \mathbf{c} and outputs a binary segmentation mask \mathbf{y} . q consists of a CNN encoder which generates a high level representation of the query image and a decoder which forms the output using skip connections, similar to the U-Net [24]. The conditional vector \mathbf{c} is fused into the query network q_{enc} at layer l through feature-wise linear modulation [16]. Afterwards, using information from the support image/segmentation, the decoder generates an output tensor of the same spatial size as the input. It consists of four blocks with channel sizes

Figure 4: Overview of MaRF (with FiLM at 3): An **encoder-decoder** network (blue) processes the query image and generates an output segmentation. Information about the search target is introduced through another **encoder** (yellow) which uses masked pooling.



of 128, 128, 128 and 32, each incorporates a skip connection from the encoder. Instead of using transposed convolutions, spatial resolution is increased using bilinear interpolation followed by a 5×5 convolution. We can describe the computation of a segmentation mask \mathbf{y} by:

$$\mathbf{y} = q_{\text{dec}}(q_{\text{enc}}^{[L:]}(q_{\text{enc}}^{[0:L]}(\mathbf{x})\phi(\mathbf{c}) + \pi(\mathbf{c})) \quad (1)$$

where the query network q_{enc} is decomposed into layers before $(q_{\text{enc}}^{[0:L]})$ and after $(q_{\text{enc}}^{[L:]})$ conditioning layer L where information from the support image is fused. Analogously to the support encoder s , the query encoder q_{enc} is implemented by a ResNet with 18 or 50 layers. Following Tian et al. [10], we replace the standard ResNet50 with a ResNet50 from PSPNet [25] in both encoder and query network. If not stated otherwise, ImageNet weights are frozen in both encoders up to the 4th residual block.

4 Automatic Output Calibration

Our model predicts a single map which expresses pixel-wise object presence probabilities. These probabilities need to be converted into binary decisions by using a threshold which separates foreground object and background. Independently from false class assignments, during our analysis, we found that the natural probability threshold of 0.5 for object presence (i.e. the separation criterion an argmax applies over two classes) is well below the optimal threshold for predictions generated by our model. However, since this underestimation of object presence happens consistently, we can compensate for it by increasing the probability for each prediction by a fixed value (or, equivalently, decreasing the object presence threshold). This results in larger fractions of the image being segmented as object (foreground). When evaluating on PASCAL-5ⁱ, instead of using a fixed threshold, we compute an individual threshold for each of the four splits of PASCAL-5ⁱ on the three unused splits respectively. Note, we apply this technique only for PASCAL-5ⁱ

5 Experiments

Datasets and Metrics We perform our analysis on three datasets: LVIS-OneShot, FSS-1000 [26] and PASCAL-

5ⁱ[7]. These datasets vary strongly in the average size and position of the object to be segmented (see Fig. 3). If not stated otherwise, both encoders q_{enc} and s_{enc} are initialized with weights obtained through ImageNet-pretraining as this is a common practice in one-shot segmentation. We exclusively use binary segmentation problems which have only a single foreground class, i.e. one-shot, one-way segmentation. We use the intersection over union-based (IoU) metrics mean IoU (mIoU), binary (or foreground-background) IoU and foreground IoU.

Training Protocol For all experiments reported in this article we used a new, slightly improved training protocol. It involves training for 120,000 iterations using a batch size of 32 with validation every 5,000 iterations. All MaRF scores reported here use FiLM conditioning at layer 3, a conditional vector size of 256 and masking at layer 4 if not stated otherwise.

LVIS-OneShot Training Instead of a fixed assignment of training image pairs, we randomly sample pairs of one category during training. Images are scaled to have a minimal side length of 480px and are then cropped to a square-shaped image. We apply mild augmentation [27] to the images involving horizontal flip as well as HSV and gamma change. For validation and test, 1000 and 10,000 fixed pairs are used without augmentation to reduce variance.

Technical Details We use PyTorch [28], Adam optimizer [29] with varying learning rates (LR), batch sizes (BS) and early stopping patience (ES) as shown below.

dataset	iter.	LR	BS
LVIS-OneShot	80,000	0.0001	32
PASCAL-5 ⁱ	5,000	0.0001	32
FSS-1000	30,000	0.0001	64

5.1 MaRF Configurations

We find that early conditioning (FiLM at layer 3) yields slightly better results on LVIS-OneShot than late conditioning (layer 5). For the smaller PASCAL-5ⁱ dataset, the opposite is true. This means a network relying on earlier conditioning is generally favorable but more data is required to

Backbone	FiLM	LVIS-OneShot		PASCAL-5 ⁱ	
		IoU _{BIN}	mIoU	IoU _{BIN}	mIoU
RN50	3	51.8	25.6	74.0	60.8
RN50	5	51.2	29.1	71.5	57.5
RN50 (no freeze)	3	58.1	34.4	71.6	57.5
RN50 (original)	3	49.2	21.9	67.9	53.5

Backbone	FiLM	IoU _{BIN}	mIoU
old:			
LVIS-OneShot	RN50	5	68.2
	RN50	3	71.4
	RN50 (no freeze)	3	67.1
	RN50 (original)	3	66.3
PASCAL-5 ⁱ	RN50	3	58.3
	RN50	5	59.4

Table 2: Performance for different configurations of MaRF trained on LVIS-OneShot. FiLM: conditioning layer L . Top new, bottom old. The new version uses LVIS training only to keep it simple. Hence the performance gap in mIoU. See link to source code in comments.

learn such a mechanism. This insight would not be visible if only a small dataset was used. Also freezing weights, like in [10], improves the performance compared to training all weights. This is likely because it prevents overfitting as LVIS-OneShot is still much smaller than image classification datasets the ResNets are normally trained on. The PSPNet [25] modification of ResNet50 turns out to be an important factor as it performs much better than conventional ResNet50 (labeled “original” in Table 2). The quality of features is an essential predictor of final performance.

5.2 Sample-efficiency and label diversity

LVIS-OneShot is larger than competitive datasets in terms of the number of categories contained. This allows us to analyze the effect of label diversity for one-shot segmentation: Given a fixed budget of samples, is there an advantage of having a *diverse* set of images? In order to answer this question we generate a set of categories C containing a specific number of samples using an iterative algorithm.

Regarding number of samples (Fig. 5, left), we find a positive relationship between number of samples and performance, which was expected as more samples generally improve performance. We observe a similarly strong correlation between performance and sample diversity (Fig. 5 right), even when the number of samples is kept constant. This result supports our intuition that not only sample size but also label diversity is crucial. The fine-grained division of classes conveys additional information useful for one-shot segmentation.

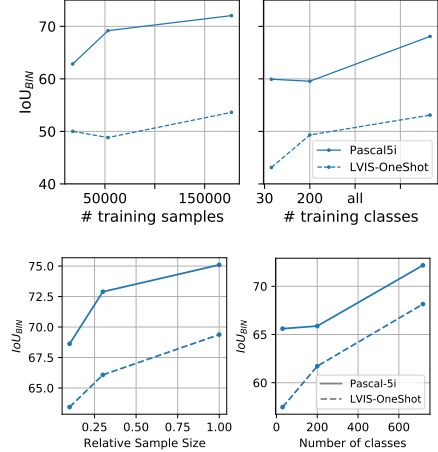


Figure 5: A larger number of samples (left) as well as more classes (right) have a positive impact on one-shot segmentation performance of MaRF (training on LVIS-OneShot). new top, old bottom

Method	Backbone	IoU _{BIN}	mIoU
Trained on PASCAL-5ⁱ			
DAN [13]	RN101	71.9	58.2
PFENet [10]	RN50	73.3	60.8
PFENet [10]	RN50	71.2*	60.6*
RPM [12]	RN50	-	56.3
RePRI [30]	RN50	-	59.7
MaRF (ours)	RN50	62.2	43.0
Trained on FSS-1000			
FSS Basel. [26]	VGG16 -	-	58.6
MaRF (ours)	RN50	67.3	45.4
Trained on LVIS-OneShot			
PFENet [10]	RN50	74.7	63.2
MaRF (ours)	RN50	74.0	60.8

Table 3: One-shot segmentation performance on PASCAL-5ⁱ. *using weights provided by the PFENet authors.

5.3 Comparison with State-of-the-art

PASCAL-5ⁱ For evaluation, we use the PASCAL-5ⁱ implementation provided by Tian et al. [10] using their training and validation sets. We find our MaRF model to perform quite poorly when it was trained on PASCAL-5ⁱ. This is expected due to its simplicity without explicit mechanisms (or inductive biases) for one-shot segmentation and the small size of the dataset.

Transfer Learning LVIS-OneShot \rightarrow PASCAL-5ⁱ In this setting we calibrate the output probabilities using the three unused splits of PASCAL-5ⁱ (there are four splits). Despite the distribution shift (Fig. 3) between the datasets, MaRF achieves competitive performance on PASCAL-5ⁱ. The quantitative results on PASCAL-5ⁱ (Table 3) indicate an

Model	10% of data		100% of data		Δ	
	IoU _{BIN}	mIoU	IoU _{BIN}	mIoU	IoU _{BIN}	mIoU
PFENet	77.5	63.7	78.0	64.1	0.5	0.4
MaRF	62.8	45.7	74.0	60.8	11.1	15.2

Table 4: Comparison between PFENet and MaRF (RN50) on the new LVIS-OneShot dataset

Model	Backbone	mIoU	mIoU _{neg}	IoU _{BIN}	IoU _{FG}
Trained on FSS-1000					
PFENet [10]	RN50 (IN)	-	-	-	80.8*
DAN [13]	RN101 (IN)	-	-	-	85.2
MaRF (FiLM 3) + aug.	RN50	81.2	42.3	88.9	83.3
MaRF (F. 3) (no support)	RN50	79.5	41.1	87.5	81.2
MaRF + aug.	RN50	70.3	37.3	80.5	71.4
MaRF	RN50 (IN)	79.8	41.5	87.1	81.0
MaRF + aug. (no supp.)	RN50 (IN)	76.3	40.3	84.7	77.3
Trained on LVIS-OneShot					
PFENet [10]	RN50	76.8	55.3	85.5	78.2
MaRF (FiLM 3, ours)	RN50	70.6	51.7	81.5	72.1
PFENet	RN50	72.0	57.6	82.4	73.2
MaRF	RN50	59.4	49.4	74.5	61.1

Table 5: Performance of MaRF (Film 3) on FSS-1000 in comparison with state-of-the-art methods. (IN) indicates ImageNet pre-training of the backbone.

mIoU score of 60.8 and IoU_{BIN} of 74.0 which is slightly better than the more complex PFENet trained on PASCAL-5ⁱ. The latter benefits from training on LVIS-OneShot, too, establishing a slightly better score. However, PFENet does not benefit to the same degree as MaRF from the larger LVIS-OneShot training dataset.

5.3.1 LVIS-OneShot

Next we investigate the relationship between performance gain and dataset size in more detail. In this setting, we use thresholds of 0.5 and 0.1 for PFENet and MaRF respectively. Consistent with our findings from the previous section, MaRF achieves a greater performance improvement from using more samples (also see Fig. 1). On one hand, this gain supports our intuition that data can outweigh model design to some extent. On the other hand, there remains a gap to PFENet, suggesting that the inductive biases of PFENet are generally useful for one-shot segmentation and are not overfit to the classes present in small segmentation datasets.

5.3.2 FSS1000

The results on FSS-1000 (Table 5) show that MaRF achieves competitive performance, only the DAN model [13] which uses a larger encoder achieves better scores. This supports the claim that simple models match state-of-the-art performance with sufficient training data, evading the need of model design.

Surprisingly, a baseline of our model that did not receive a support image/segmentation (labeled “no supp.” in Ta-

ble 5) achieved decent performance, and even outperformed all previously published approaches except DAN [13]. This result suggests that FSS-1000 is strongly biased towards centered objects and has little variation in object size and location.

To further investigate the biases of FSS-1000, we introduce 50% negative samples to the test set. We observe a strong drop in performance, while the models trained on LVIS-OneShot cope best with this setting. Possibly due to these different statistics, transfer learning from LVIS-OneShot does not work as well for FSS-1000 as for PASCAL-5ⁱ.

6 Discussion and Conclusion

Large-scale training can replace model design and strong inductive biases in one-shot segmentation. This result is consistent with previous findings in computer vision [31, 4, 6] and NLP [32, 33]. We find conceptually simple models to profit to a much greater extent from more samples and more diverse samples than the complex PFENet. The latter achieved only small gains from a substantial increase of dataset richness. However, PFENet performance remains slightly better than our baseline, suggesting that inductive biases still matter in the large data-regime (although to a smaller extent) and PFENet modeled the right ones. For future research in one-shot semantic segmentation, our findings represent a strong argument in favor of using large and diverse datasets. We recommend to consider PASCAL-5ⁱ primarily a test dataset.

References

- [1] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems (NIPS)*, 2012.
- [2] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2009.
- [3] Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly, and Neil Houlsby. Big transfer (bit): General visual representation learning. *arXiv preprint arXiv:1912.11370*, 6, 2019.
- [4] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020*, 2021.
- [5] Yonglong Tian, Yue Wang, Dilip Krishnan, Joshua B. Tenenbaum, and Phillip Isola. Rethinking few-shot image classifi-

- cation: A good embedding is all you need? In *European Conference on Computer Vision (ECCV)*, Cham, 2020. Springer International Publishing.
- [6] Claudio Michaelis, Matthias Bethge, and Alexander S. Ecker. Closing the generalization gap in one-shot object detection. 2020.
 - [7] Amirreza Shaban, Shray Bansal, Zhen Liu, Irfan Essa, and Byron Boots. One-shot learning for semantic segmentation. *BMVC*, 2017.
 - [8] Khoi Nguyen and Sinisa Todorovic. Feature weighting and boosting for few-shot segmentation. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019.
 - [9] Qiang Wang, Li Zhang, Luca Bertinetto, Weiming Hu, and Philip HS Torr. Fast online object tracking and segmentation: A unifying approach. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1328–1338, 2019.
 - [10] 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 and machine intelligence*, PP, August 2020. ISSN 0162-8828. doi: 10.1109/tpami.2020.3013717. URL <https://doi.org/10.1109/TPAMI.2020.3013717>.
 - [11] Yongfei Liu, Xiangyi Zhang, Songyang Zhang, and Xuming He. Part-aware prototype network for few-shot semantic segmentation. In *European Conference on Computer Vision*, pages 142–158. Springer, 2020.
 - [12] 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 (ECCV)*, 2020.
 - [13] Xianbin Cao and Xiantong Zhen. Few-shot semantic segmentation with democratic attention networks. *European Conference on Computer Vision (ECCV)*, 2020.
 - [14] Agrim Gupta, Piotr Dollar, and Ross Girshick. Lvis: A dataset for large vocabulary instance segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
 - [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
 - [16] Vincent Dumoulin, Ethan Perez, Nathan Schucher, Florian Strub, Harm de Vries, Aaron Courville, and Yoshua Bengio. Feature-wise transformations. *Distill*, 2018. doi: 10.23915/distill.00011.
 - [17] Chi Zhang, Guosheng Lin, Fayao Liu, Jiushuang Guo, Qingyao Wu, and Rui Yao. Pyramid graph networks with connection attentions for region-based one-shot semantic segmentation. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9587–9595, 2019.
 - [18] Reza Azad, Abdur R Fayjie, Claude Kauffmann, Ismail Ben Ayed, Marco Pedersoli, and Jose Dolz. On the texture bias for few-shot cnn segmentation. In *IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, 2021.
 - [19] Kaixin Wang, Jun Hao Liew, Yingtian Zou, Daquan Zhou, and Jiashi Feng. Panet: Few-shot image semantic segmentation with prototype alignment. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9197–9206, 2019.
 - [20] Khoi Nguyen and Sinisa Todorovic. Feature weighting and boosting for few-shot segmentation. In *IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 622–631, 2019.
 - [21] 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 (ECCV)*, 2014. ISBN 978-3-319-10601-4.
 - [22] George A Miller. Wordnet: a lexical database for english. 38 (11):39–41, 1995.
 - [23] Xiaolin Zhang, Yunchao Wei, Yi Yang, and Thomas S Huang. Sg-one: Similarity guidance network for one-shot semantic segmentation. *IEEE Transactions on Cybernetics*, 2020.
 - [24] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical Image Computing and Computer-assisted Intervention*. Springer, 2015.
 - [25] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
 - [26] Tianhan Wei, Xiang Li, Yau Pun Chen, Yu-Wing Tai, and Chi-Keung Tang. Fss-1000: A 1000-class dataset for few-shot segmentation. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
 - [27] Alexander Buslaev, Vladimir I Iglovikov, Eugene Khvedchenya, Alex Parinov, Mikhail Druzhinin, and Alexandr A Kalinin. Albumentations: fast and flexible image augmentations. *Information*, 2020.
 - [28] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.

- [29] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [30] 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? *arXiv preprint arXiv:2012.06166*, 2020.
- [31] Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. Revisiting unreasonable effectiveness of data in deep learning era. In *IEEE International Conference on Computer Vision (ICCV)*, pages 843–852, 2017.
- [32] Alon Halevy, Peter Norvig, and Fernando Pereira. The unreasonable effectiveness of data. *IEEE Intelligent Systems*, 24(2):8–12, 2009.
- [33] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI blog*, 2019.
- [34] Timo Luddecke and Alexander Ecker. The role of data for one-shot semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pages 2653–2658, June 2021.
- [35] Steven Bird. Nltk: the natural language toolkit. In *Proceedings of the COLING/ACL 2006 Interactive Presentation Sessions*, pages 69–72, 2006.

Detailed Correction

Background

In a previous version of this paper [34] we introduced the LVIS-OneShot dataset involving a larger number of classes than previous one-shot segmentation datasets and found the conceptually simple MaRF model trained on this dataset to perform well in the Pascal-5i benchmark. However, we discovered a problem in the assignment of object classes to training and test splits of LVIS-OneShot, which should contain Pascal classes only in the test split. As a result of this problem, the following classes were falsely assigned to the training dataset: *television set, jet plane, fighter jet, motorcycle, airplane, person, cow, dirt bike*. As a consequence, since objects of these classes were seen during training, the reported scores computed on Pascal-5i do not measure one-shot performance.

Error Analysis We implemented an automatic matching procedure based on WordNet (NLTK implementation [35]) which traverses hypernyms (i.e. more broader concepts) of all LVIS object classes. If a match between the hypernyms of an LVIS object class and a Pascal class is found, the LVIS object class is assigned to the test set. For this we manually assigned all 20 pascal classes to Wordnet synsets (a synset corresponds to one semantic meaning). Here we mistakenly selected unfavorable synsets for which the hypernym search failed. For example, for we used *cow.n.01* as the synset for cows in Pascal, but in LVIS cows are associated with the synset *beef.n.01*. Since both, cow and beef are hyponyms (i.e. more special categories) of cattle, traversing the hypernyms of *beef.n.01* did not produce a match with *cow.n.01*. Similar problems led to the false assignment of the other synsets.

Correction In addition to the synset assignment problem we found that the output probabilities of our model are often far from the optimum. By choosing lower thresholds than 0.5, we are able to obtain much better results. In general a threshold of 0.1 works well. For PASCAL-5ⁱ we compute optimal thresholds for each of its four splits using the remaining three training splits. We fixed the synset assignment algorithm, implemented the output probability calibration and re-ran the affected experiments using the new training protocol. As a consequence of correcting the training split, we observed a substantial drop in performance by over 10 percent points. However, much of this decrease is compensated by the positive effect of the calibration. Note, this calibration was not conducted in the original paper [34] (it would have improved scores of MaRF to around 64%). We further investigate two questions regarding calibration: (1) *Do we need calibration because of randomly sampling training pairs?* Instead of using a fixed training set, we ran-

domly draw query image and support image (of the same category) during training. In order to investigate the effect of this random sampling, we use a checkpoint of our model and continue training on a set of 50.000 fixed query-support image pairs for training. While the training loss decreased much faster than with random pair sampling to around 0.08 (compared to around 0.14 with random sampling), general performance remained slightly lower than with random sampling. Importantly, we found that probabilities still need to be corrected. (2) *Is calibration needed for all models?* We also applied the calibration on the weights provided by the PFENet authors. Here we found that calibration did not help, indeed it slightly decreased performance.

We conclude that the necessity for calibration comes mainly from unknown classes. Since these classes have not been seen during training, the model is less confident and thus predicts smaller logits. In PFENet, the high-level feature similarity prior likely compensates for this.

Conclusion

After removing the Pascal classes from the training set, performance scores drop dramatically. However, this gap can be compensated to a large extent by calibrating the output probabilities such that a larger fraction of the foreground is segmented. The finding that not only more samples but also more class diversity leads to better performance remains true. The simple MaRF model trained on LVIS-OneShot does no longer outperform the more complex PFENet model trained on Pascal-5i but the gap is fairly small. Our general recommendation of using larger and more diverse datasets than Pascal-5i for one-shot segmentation remains intact.

	MaRF			PFENet [10]		
	mIoU	IoU _{FG}	IoU _{BIN}	mIoU	IoU _{FG}	IoU _{BIN}
reported in [34] (wrong splits)	61.1	-	77.4	64.3	-	78.5
corrected splits	40.1	41.7	64.4	56.6	58.7	74.3
corrected splits + calibration	60.2	59.2	73.4	63.4	63.8	76.7
trained on Pascal-5i	38.3	37.6	60.9	60.0*	55.3*	70.7*
trained on Pascal-5i + calibration	42.0	40.5	62.2	58.6*	54.8*	70.8*

Table 6: Pascal-5i mIoU performance in %. Values indicated by * were obtained by us using the weights provided by the PFENet authors. [TODO: this old table is just kept for reference and will be removed](#)