

# A Strong Baseline for Generalized Few-Shot Semantic Segmentation

Sina Hajimiri\*

Malik Boudiaf

Ismail Ben Ayed

Jose Dolz

ÉTS Montreal

## Abstract

*This paper introduces a generalized few-shot segmentation framework with a straightforward training process and an easy-to-optimize inference phase. In particular, we propose a simple yet effective model based on the well-known InfoMax principle, where the Mutual Information (MI) between the learned feature representations and their corresponding predictions is maximized. In addition, the terms derived from our MI-based formulation are coupled with a knowledge distillation term to retain the knowledge on base classes. With a simple training process, our inference model can be applied on top of any segmentation network trained on base classes. The proposed inference yields substantial improvements on the popular few-shot segmentation benchmarks, PASCAL-5<sup>i</sup> and COCO-20<sup>i</sup>. Particularly, for novel classes, the improvement gains range from 7% to 26% (PASCAL-5<sup>i</sup>) and from 3% to 12% (COCO-20<sup>i</sup>) in the 1-shot and 5-shot scenarios, respectively. Furthermore, we propose a more challenging setting, where performance gaps are further exacerbated. Our code is publicly available at <https://github.com/sinahmr/DIaM>.*

## 1. Introduction

With the advent of deep learning methods, the automatic interpretation and semantic understanding of image content have drastically improved in recent years. These models are nowadays at the core of a broad span of visual recognition tasks and have enormous potential in strategic domains for our society, such as autonomous driving, healthcare, or security. Particularly, semantic segmentation, whose goal is to assign pixel-level categories, lies as one of the mainstays in visual interpretation. Nevertheless, the remarkable performance achieved by deep learning segmentation models is typically limited by the amount of available training data. Indeed, standard segmentation approaches are often trained on a fixed set of predefined semantic categories, commonly requiring hundreds of examples per class. This limits their scalability to novel classes, as obtaining annotations for new categories is a cumbersome and labor-intensive process.

Few-shot semantic segmentation (FSS) has recently emerged as an appealing alternative to overcome this limitation [1, 30, 33]. Under this learning paradigm, models are trained with an abundant labeled dataset on *base* classes, and only a few instances of *novel* classes are seen during the adaptation stage. However, [29] identified two important limitations that hamper the application of these methods in real-life scenarios. First, existing literature on FSS assumes that the support samples contain the categories present in the query images, which may incur costly manual selection processes. Second, even though significant achievements have been made, all these methods focus on leveraging supports as much as possible to extract effective target information, but neglect to preserve the performance on known categories. Furthermore, while in many practical applications the number of novel classes is not limited, most FSS approaches are designed to work on a binary basis, which is suboptimal in the case of multiple novel categories.

Inspired by these limitations, a novel Generalized Few-Shot Semantic Segmentation (GFSS) setting has been recently introduced in [29]. In particular, GFSS relaxes the strong assumption that the support and query categories are the same. This means that, under this new learning paradigm, providing support images that contain the same target categories as the query images is not required. Furthermore, the evaluation in this setting involves not only *novel* classes but also *base* categories, which provides a more realistic scenario.

Although the setting in [29] overcomes the limitations of few-shot semantic segmentation, we argue that a gap still remains between current experimental protocols and real-world applications. Hereafter, we highlight limiting points of the current literature and further discuss them in Sec. 3.2.

**Unrealistic prior knowledge.** We found that existing works explicitly rely on prior knowledge of the novel classes (supposed to be seen at test-time only) during the training phase. This, for instance, allows to filter out images containing novel objects [14, 29] from the training set. Recent empirical evidence [28] found out that such assumptions indeed boost the results in a significant manner.

\*Corresponding author: seyed-mohammadsina.hajimiri.1@etsmtl.net

**Modularity.** Another limitation is the tight entanglement between the training and testing phases of current approaches, which often limits their ability to handle arbitrary tasks at test time. Specifically, existing meta-learning-based approaches are designed to handle binary segmentation [14], and need to be consequently modified to handle multiple classes. While we technically address that by using multiple forward passes (one per class) followed by some heuristic aggregation of segmentation maps, this scales poorly and lacks principle.

**Contributions.** Motivated by these limitations, we aim to address a more practical setting and develop a fully modular inference procedure. Our inference abstracts away the training stage, making no assumption about the type of training or the format of tasks met at test time. Specifically:

- We present a new GFSS framework, *DIaM* (Distilled Information Maximization). Our method is inspired by the well-known InfoMax principle, which maximizes the Mutual Information between the learned feature representations and their corresponding predictions. To reduce performance degradation on the base categories, without requiring explicit supervision, we introduce a Kullback-Leibler term that enforces consistency between the old and new model’s base class predictions.
- Although disadvantaged by rectifications to improve the practicality of previous experimental protocols, we still demonstrate that *DIaM* outperforms current SOTA on existing GFSS benchmarks, particularly excelling in the segmentation of novel classes.
- Based on our observations, we go beyond standard benchmarks and present a more challenging scenario, where the number of base and novel classes is the same. In this setting, the gap between our method and the current GFSS SOTA widens, highlighting the poor ability of modern GFSS SOTA to handle numerous novel classes and the need for more modular/scalable methods.

## 2. Related work

**Few-shot segmentation.** Few-shot semantic segmentation (FSS) has received notable attention in recent years, greatly inspired by the success of the few-shot learning paradigm [9, 24]. Early FSS frameworks consisted in a dual-branch architecture, where one branch generated the class prototypes from support samples and the other one segmented the query images by exploiting the learned prototypes [7, 23, 25]. Following the success of these pioneer approaches, an important body of literature explored how to better leverage category information from support

samples to better guide the segmentation of query images [15, 21, 30, 33, 34, 36]. For example, this can be achieved by collecting more abundant information from support images, which is used to construct multiple prototypes per class, each activating different regions of the query image [15, 34]. Alternative solutions to learn better category representations include: establishing correspondences between support and query images with Graph CNNs [32], imprinting the weights for novel classes [27], or leveraging visual transformers to improve the category information transfer between support and query samples [20, 22, 38]. Last, inspired by recent works in few-shot classification that favor a transductive setting, foregoing episodic training (*aka* meta-learning) [2, 5, 19, 40], RePRI [1] proposed a simple transductive solution.

**Generalized few-shot segmentation.** To overcome some of the limitations of FSS, [29] recently extended this setting, which was coined as generalized few-shot semantic segmentation (GFSS). In particular, GFSS approaches are given a single support set containing some images for every novel class, and they should be able to predict all potential base and novel classes in all query images. This way, in contrast to standard FSS methods, models have no knowledge of novel classes present in a query image. To tackle this problem, CAPL [29] proposed a framework with two modules to dynamically adapt both base and novel prototypes. Nevertheless, the presented results are biased toward base classes and the solution requires that base classes are labeled in the support samples. Furthermore, the recent BAM model [14], which was initially proposed for FSS, is also evaluated in the GFSS setting. This model consists of two steps. First, a *base-learner* is trained on base classes following the standard supervised learning paradigm, where the cross entropy loss is employed on the base training set. Then, a second meta-learning step is introduced, where the base-learner and a new *meta-learner* are optimized using episodic training. In the inference phase, the output of the meta-learner is fused with base-learner’s output to give predictions on base classes and a single novel class. The fact that the meta-learner is only able to discern background-foreground categories makes this method’s direct application not suitable to multi-class GFSS.

## 3. Background

### 3.1. Preliminaries

**Notations.** Let us note  $H$  and  $W$  the height and width of the original images, and  $\Omega = [0, H - 1] \times [0, W - 1]$  the set of all pixels coordinates. In all generality, we define a segmentation model  $f$  that takes images  $\mathbf{x} \in \mathbb{R}^{|\Omega| \times 3}$  as inputs, and produces segmentation maps  $f(\mathbf{x}) = \mathbf{p} \in [0, 1]^{|\Omega| \times K}$ , where  $K$  denotes the number of classes to predict.

**Standard few-shot segmentation.** In few-shot segmentation, two sets of classes are considered: the *base* classes,  $\mathcal{C}^b$ , containing classes over which the model is trained; and the *novel* classes,  $\mathcal{C}^n$ , strictly disjoint from base classes, such that  $\mathcal{C}^b \cap \mathcal{C}^n = \emptyset$ . At test time, the model is evaluated through a series of tasks. In each task, the model is given access to a support set  $\mathbb{S} = \{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^{|\mathbb{S}|}$ , containing a few images (shots), along with their corresponding binary segmentation masks  $\mathbf{y}_i \in [0, 1]^{|\Omega|}$ , for some novel class randomly sampled from  $\mathcal{C}^n$ . Using this limited supervision, the model is then evaluated based on its ability to segment objects from this novel category in an unlabeled image, referred to as the *query* image, which will be referred to as the  $(|\mathbb{S}| + 1)^{th}$  image. In this context,  $K = 2$  and the model is expected to produce a binary mask  $\mathbf{p} \in [0, 1]^{|\Omega| \times 2}$ .

**Generalized few-shot segmentation.** The *generalized* setting extends the standard setting to account for the fact that real-world applications often require being able to recognize both base and novel classes in new images. As a result, for each task, we now require the model to produce a segmentation over the  $1 + |\mathcal{C}^b| + |\mathcal{C}^n|$  (including the background) potential classes, such that  $\mathbf{p} \in [0, 1]^{|\Omega| \times (1 + |\mathcal{C}^b| + |\mathcal{C}^n|)}$ . So, for a given pixel  $j$  we have

$$\mathbf{p}(j) = \left[ \underbrace{p_0}_{bg}, \underbrace{p_1, \dots, p_{|\mathcal{C}^b|}}_{\text{base classes}}, \underbrace{p_{|\mathcal{C}^b|+1}, \dots, p_{|\mathcal{C}^b|+|\mathcal{C}^n|}}_{\text{novel classes}} \right]^\top, \quad (1)$$

in which we omit pixel index  $j$  from the right-hand side for simplicity, and *bg* stands for the background.

### 3.2. Toward a fully practical setting

As motivated in Sec. 1, we aim at evaluating methods in a maximally practical setting. Therefore, we start by rectifying design choices made in previous works [14, 29] that we find impractical and that could impact the results in a significant manner. We further posit additional desiderata to improve the practicality of developed methods and widen their adoption.

**Addressing the presence of novel classes during training.** We found that previous works [14, 29] explicitly removed images containing novel classes during the training phase. This implicitly requires information that should not be available at that stage, namely the prior knowledge of the novel classes, as well as the potential presence of a given novel object in a particular image. *Instead, in our setting, we keep those images during training (as they naturally occur), and the potential objects from novel classes are labeled as background at that stage.* Needless to say, this may negatively impact the performance of the model at test time, given the class ambiguity introduced while forcing the

network to predict potential novel classes as background. However, we believe that this is a more natural way to design the problem.

**Relaxing test-time labeling requirements.** We found that previous works [29] required potential objects from base classes to be explicitly labeled in the images of the support set  $\mathbb{S}$ . We argue that this can require a significant additional load of work in real-world settings. Consider the simple example of COCO-20<sup>i</sup>'s dataset, with a total of 80 classes. Instead of only annotating objects from 20 novel classes, as was the case in the standard FSS setting, annotators would now have to search for 60 potential additional classes in each image from the support, thus leading to a substantial increase in human/financial efforts. We find that this requirement is not necessary, and high performances on base classes can be retained by other means. *Therefore, we drop this requirement, and only require annotations for novel classes at test-time, whereas the rest (including potential objects from base classes) are labeled as background.*

**Modularity of inference.** Although not a strict requirement, we advocate developing modular inferences that, unlike current approaches [14, 29], can apply to any model, without relying on customized architectures or training procedures. The rationale is two-fold. First, as foundation models are and will continue to push state-of-the-art on most vision tasks, we forecast that the ability to leverage off-the-shelf models seamlessly will become crucial in reaching high performances. Second, it drastically lowers the entry barrier to few-shot learning for practitioners who, in most cases, already possess trained models for their specific application, and whose limited computational resources may prevent re-training models for every method they would like to try. *Access to inference-only methods which can readily equip preexisting models with few-shot ability could be a key ingredient for the widespread adoption of few-shot methods.*

## 4. Our method

In light of the requirements and desiderata posed in Sec. 3.2, we shift our attention from training to inference. In particular, unlike previous generalized FSS methods, we use a standard supervised training procedure that is not informed by any knowledge, even implicit, of the novel classes. We further develop an optimization-based inference procedure that can be directly deployed at test time. Each of these steps is detailed below.

### 4.1. Training

For convenience, we partition the segmentation model into a feature extractor  $f_\phi$  and a linear classifier  $f_{\theta_\phi}$ , trained

in a standard supervised fashion to segment base classes  $\mathcal{C}^b$  during training. At this stage, the classifier can only predict  $1 + |\mathcal{C}^b|$  classes, *i.e.*, the background and the base classes.

## 4.2. Inference

At test-time, given  $|\mathcal{C}^n|$  novel classes to recognize, we freeze the feature extractor  $f_\phi$  and augment the pre-trained classifier  $\theta_b \in \mathbb{R}^{(1+|\mathcal{C}^b|) \times d}$  with novel prototypes  $\theta_n \in \mathbb{R}^{|\mathcal{C}^n| \times d}$ . We consider the concatenation  $\theta = [\theta_b; \theta_n] \in \mathbb{R}^{(1+|\mathcal{C}^b|+|\mathcal{C}^n|) \times d}$  to form our final classifier, and optimize  $\theta$  for this specific task. Note that  $d$  is the size of the feature space. We base our optimization objective on the seminal concept of mutual information [26]. Specifically, we use the InfoMax framework [18] as a starting point for our formulation. InfoMax advocates maximizing the mutual information between a network’s inputs and outputs as

$$\max_{\theta} I(X; P) = \underbrace{H(P)}_{\text{marginal entropy}} - \underbrace{H(P|X)}_{\text{conditional entropy}}, \quad (2)$$

where  $X$  and  $P$  are the random variables respectively associated with the pixel distribution and model’s predictions. Instantiated in our context, InfoMax (2) incites the model to produce confident predictions on each pixel (conditional entropy), while encouraging an overall balanced marginal distribution (marginal entropy), *i.e.*, roughly speaking, an equal number of pixels assigned to each class. Interestingly, in the related context of classification, InfoMax can be interpreted as an unsupervised clustering criterion [13].

In the following sections, we explain how we gradually depart from the vanilla InfoMax principle and incorporate problem-specific constraints and inductive biases to reach our final formulation.

### 4.2.1 Enforcing high-confidence

We start by focusing our attention on the conditional entropy term mentioned in Eq. (2) that enforces high-confidence predictions for each pixel. To this end, we introduce the cross-entropic operator:

$$H(\mathbf{p}_i; \mathbf{q}_i) = \frac{-1}{|\Omega|} \text{Tr}(\mathbf{p}_i \log(\mathbf{q}_i^\top)). \quad (3)$$

Taking into account support and query images, the vanilla conditional entropy of  $P|X$  can be written as

$$H(P|X) = \frac{1}{|\mathcal{S}| + 1} \sum_{i=1}^{|\mathcal{S}|+1} H(\mathbf{p}_i; \mathbf{p}_i). \quad (4)$$

**Leveraging supervision** Entropy  $H(\mathbf{p})$  can be interpreted as a self-cross-entropy  $H(\mathbf{p}; \mathbf{p})$ , in which the model’s own predictions are used as pseudo-labels for supervision. Because actual ground-truth labels for the support images are

provided, we can effectively replace those pseudo-labels with the ground truth.

**Aligning support labels and predictions.** As motivated in Sec. 3.2, we do not require objects from base classes to be labeled in the support images. That produces a slight misalignment between labels and the model’s predictions. More specifically, a pixel  $j$  labeled as *background*,  $\mathbf{y}_i(j) = [1, 0, \dots, 0]$ , can now have two meanings: either it actually is a background pixel or it belongs to a base class object. To account for that misalignment between predictions  $\mathbf{p}$  and labels  $\mathbf{y}$ , we project the model’s predictions as

$$\pi_{\mathcal{S}}(\mathbf{p}_i)(j) = \left[ \sum_{k=0}^{|\mathcal{C}^b|} p_k, \underbrace{0, \dots, 0}_{|\mathcal{C}^b| \text{ times}}, p_{|\mathcal{C}^b|+1}, \dots, p_{|\mathcal{C}^b|+|\mathcal{C}^n|} \right]^\top. \quad (5)$$

We can now write our constrained conditional entropy, partitioning pixels into supervised and unsupervised. We found it beneficial to adjust the relative weighting of the terms, and therefore introduce  $\alpha > 0$  in the objective, such that our conditional entropy term reads as

$$\mathcal{L}_{\text{cond-ent}} = \alpha \underbrace{\sum_{i=1}^{|\mathcal{S}|} H(\mathbf{y}_i; \pi_{\mathcal{S}}(\mathbf{p}_i))}_{\mathcal{L}_{\text{xent}}: \text{Support supervised entropy}} + \underbrace{H(\mathbf{p}_{|\mathcal{S}|+1})}_{\text{Query unsupervised entropy}}, \quad (6)$$

where  $\alpha$  controls the reliance on the labeled support set. We will refer to the *support supervised entropy* term as  $\mathcal{L}_{\text{xent}}$ .

### 4.2.2 Addressing class imbalance

The constraint of high-confidence predictions alone can be easily satisfied by a model and does not constrain the problem enough to guarantee meaningful solutions. For instance, a trivial classifier assigning all pixels from the query image to the same class with maximum probability 1 would fully satisfy the high-confidence constraint. Therefore, we need to go beyond mere entropy minimization, which naturally leads us to shift our attention to the marginal entropy term from Eq. (2). As shown below, marginal entropy ensures a fair distribution of assignments over the different classes, thereby preventing, *e.g.*, the trivial solutions previously described. As in the previous section, let us start with the vanilla formulation of the marginal entropy

$$H(P) = -\hat{\mathbf{p}} \cdot \log(\hat{\mathbf{p}}) = \text{Cste} - \text{KL}(\hat{\mathbf{p}} \parallel \mathbf{u}), \quad (7)$$

where  $\cdot$  is the dot-product,  $\text{KL}(\cdot \parallel \cdot)$  denotes the Kullback-Leibler divergence,  $\mathbf{u} = 1/(1 + |\mathcal{C}^b| + |\mathcal{C}^n|) \cdot \mathbf{1}$  is the uniform distribution, and  $\hat{\mathbf{p}} \in [0, 1]^{1+|\mathcal{C}^b|+|\mathcal{C}^n|}$  is



the model’s marginal distribution over classes (detailed hereafter). In other words,  $H(P)$  regularizes the overall procedure by encouraging class-balanced predictions, *i.e.*, predicting roughly an even distribution of pixels for each class. In the context of our problem, this vanilla formulation exhibits important limitations.

First, akin to the previous section, we have direct access to supervision for the support pixels. Assuming that the model fits those labels well, the distribution of predictions will naturally converge to the proportions dictated by the ground truth labels. Therefore, we do not include the predictions from the support samples when computing the marginal, as this would provide a redundant signal, and only consider the marginal distribution over the query:

$$\hat{\mathbf{p}} = \frac{1}{|\Omega|} \sum_{j \in \Omega} \mathbf{p}_{|\mathcal{S}|+1}(j). \quad (8)$$

Additionally, semantic segmentation is virtually never a balanced problem. The number of instances, the distance to the camera, or the angle of view are all factors that can randomly vary between scenes, significantly affecting the final share of pixels that each class occupies in a given frame. Therefore, using the uniform distribution as a prior to match can be sub-optimal, as shown in our ablation study in Sec. 5.3. Beyond merely down-weighting to weaken this regularization, alternatives in the literature include using  $\alpha$ -entropy [31] in place of the standard Shannon entropy from Eq. (7), or replacing  $\mathbf{u}$  with an estimated prior  $\mathbf{\Pi}$  [1]. We decided to go with the prior estimation procedure given in [1]. Specifically, we extend the marginal entropy to take into account a prior:

$$\mathcal{L}_{\text{marg-ent}} = H(P; \mathbf{\Pi}) = \text{Cste} - \text{KL}(\hat{\mathbf{p}} \parallel \mathbf{\Pi}). \quad (9)$$

This *prior-guided* marginal entropy loss reduces to the standard marginal entropy in the absence of prior, *i.e.*,  $\mathbf{\Pi} = \mathbf{u}$ . Following [1],  $\mathbf{\Pi}$  is estimated from the model’s initial marginal distribution and re-updated during optimization.

#### 4.2.3 Preserving base knowledge

So far, our inference procedure has not made any distinction between base and novel classes, thus leaving aside an important inductive bias of our problem: prototypes from the base classifier,  $\theta_b$ , were trained using orders of magnitude more data than prototypes from novel classes, whose only available supervision comes from the few labeled samples from the support set. Additionally, our setting prevents support images from providing any explicit supervision for base classes, which only accentuates the asymmetry between base and novel classes.

To account for these two contrasts, a simple solution could be to freeze the base classifier  $\theta_b$ , and only optimize the novel classifier  $\theta_n$ . However, we show in Sec. 5 that this results in sub-optimal results. Instead, we propose a more flexible self-distillation term that encourages the model’s predictions on base classes to stay close to its old predictions. Formally, we consider the base classifier’s weights  $\theta_b^{(0)}$  (right after training) and define the model’s *old* predictions as

$$\mathbf{p}_i^{\text{old}} = f_{\theta_b^{(0)}} \circ f_{\phi}(\mathbf{x}_i) \in [0, 1]^{|\Omega| \times (1+|\mathcal{C}^b|)}. \quad (10)$$

**New-to-old mapping.** In order to measure and minimize any sort of distances between the old model’s predictions  $\mathbf{p}_i^{\text{old}}$  defined over base classes and our current model’s predictions  $\mathbf{p}_i \in [0, 1]^{|\Omega| \times (1+|\mathcal{C}^b|+|\mathcal{C}^n|)}$ , defined over both base and novel classes, we must map them to the same label space. At this point, it is important to recall that a *background* prediction from the base model really means “anything other than base classes”. That includes actual background, as well as potential novel classes that were labeled as background during training. Therefore, to make  $\mathbf{p}$  and  $\mathbf{p}^{\text{old}}$  consistent, we project  $\mathbf{p}$  as

$$\pi_{\text{new2old}}(\mathbf{p})(j) = \left[ p_0 + \sum_{i=1}^{|\mathcal{C}^n|} p_{|\mathcal{C}^b|+i}, p_1, p_2, \dots, p_{|\mathcal{C}^b|} \right]^{\top}. \quad (11)$$

Now, inspired by recent literature in incremental learning that distills knowledge using the predictions of old models [3, 4, 6, 16], we can express our *knowledge-distillation* term applied to the query image as

$$\mathcal{L}_{\text{KD}} = \text{KL}(\pi_{\text{new2old}}(\mathbf{p}_{|\mathcal{S}|+1}) \parallel \mathbf{p}_{|\mathcal{S}|+1}^{\text{old}}), \quad (12)$$

which enables us to write our final objective:

$$\min_{\theta} \mathcal{L}_{\text{DlaM}} = \mathcal{L}_{\text{cond-ent}} - \mathcal{L}_{\text{marg-ent}} + \beta \mathcal{L}_{\text{KD}}, \quad (13)$$

where  $\beta$  controls the importance of retaining base classes knowledge.

## 5. Experiments

### 5.1. Experimental setting

**Datasets.** To evaluate our method we use two well-known few-shot segmentation benchmarks: PASCAL-5<sup>i</sup> [8, 10, 25] and COCO-20<sup>i</sup> [17, 25]. To use in our experiments in Tab. 2, we define PASCAL-10<sup>i</sup> in the same way PASCAL-5<sup>i</sup> is formed [25], but splitting the set of classes into two subsets of size 10, instead of four subsets of size 5. More specifically, in PASCAL-10<sup>i</sup>, the subset  $i$  consists of classes with indices  $\{10i + j\}$  for  $j \in \{1, 2, \dots, 10\}$ . For COCO-20<sup>i</sup>

we report the average performance of models over 10K query images, while for PASCAL-5<sup>i</sup> and PASCAL-10<sup>i</sup> we use all available query images. More specifications about the datasets can be found in [Appendix A](#).

**Evaluation protocol.** For evaluation purposes, we resort to the standard mean intersection-over-union (mIoU) over the classes. In our tables, *Base* and *Novel* refer to mIoU over base and novel classes, respectively. Although mIoU over all classes has been used by prior works [14, 29] as the overall score, we believe it is a misleading metric in GFSS. Classes in PASCAL-5<sup>i</sup> and COCO-20<sup>i</sup> are split in a way that the number of base classes is thrice the number of novel classes. This biases the metric toward the *Base* score and undermines the very goal of few-shot learning, which is learning the novel classes. Therefore, we propose to use the standard average of *Base* and *Novel* scores as the *Mean* score in the GFSS task. Following [14, 29], metrics are averaged over 5 independent runs.

**Implementation details.** The architecture of the model is based on PSPNet [39] using Resnet-50 [11] backbone. During training, a standard cross-entropy over the base classes is minimized. Our training scheme follows the base-training stage of [14], so, the batch size is 12 and SGD optimizer is used with an initial learning rate  $2.5 \times 10^{-4}$ , momentum 0.9, and weight decay  $10^{-4}$ . The number of epochs is 20 for COCO-20<sup>i</sup> and 100 for PASCAL-5<sup>i</sup> and PASCAL-10<sup>i</sup>. Data augmentation is done in the same way as [30]. At inference time, the feature extractor  $\phi$  is kept frozen and the classifier  $\theta$  is optimized. For this phase, SGD optimizer is used with learning rate  $1.25 \times 10^{-3}$  and the loss function in Eq. (13) is optimized for 100 iterations. The size of the feature space  $d$  is set to 512 in all experiments. We have empirically found that  $\mathcal{L}_{\text{xent}}$  in Eq. (6) and  $\mathcal{L}_{\text{KD}}$  of Eq. (12) play more significant roles in the model’s performance, and upweighted these terms by two orders of magnitude ( $\alpha = \beta = 100$ ). Following [1], the value of  $\Pi$  is estimated by the model at the beginning of the evaluation and it is updated once at iteration 10.

**Baselines.** Following [29], we included relevant FSS methods in our evaluation, including CANet [37], PANet [33], PFENet [30], and SCL [36]. These methods were adapted in [29] by modifying their respective inference code to generate prototypes for both base and novel classes in each query image. We also modified the inference code of RePRI [1] to accommodate multiple classes during testing and adapted MiB [4], an incremental learning method, to the GFSS setting, as a distillation term similar to Eq. (12) is integrated in their approach. Furthermore, we compare our method to the GFSS method CAPL [29]. Last, although BAM [14] reports results in the GFSS task, its

episodic learning nature hinders the scalability of this approach to settings where segmentation of multiple novel classes is required. Indeed, at inference, it can only provide background-foreground predictions, which is impractical in our current validation. In order to include it in our experiments, we have made some changes to this method, which are detailed in [Appendix G](#). Note that the reported results in Tab. 1 do not take into account the background IoU in the evaluation metrics, as this class is not an object of interest. This is discussed in more detail in [Appendix E](#).

## 5.2. Main results

**Comparison to state-of-the-art GFSS.** Table 1 reports the results obtained by different approaches in the GFSS setting. Here, we stress some of the underlying limitations present in these methods. First, we refer to as ‘Practical setting’ the scenario where models employ the whole dataset during the training of base classes, *i.e.*, have no access beforehand to information about novel categories (more information is available in [Appendix F](#)). Then, we resort to ‘Multi-class design’ to highlight which methods, by nature, can handle multiple novel classes simultaneously. From these results, we can observe that models designed specifically for the task of FSS are unable to perform satisfactorily in both base and novel categories. Compared to GFSS methods, our formulation brings substantial improvements under both the 1-shot and 5-shot scenarios, particularly on *Novel* metric. More specifically, compared to CAPL, these differences are considerably large in the case of novel classes, with around 17% and 31% improvement on PASCAL-5<sup>i</sup>. We believe that this imbalanced behavior might be due to the fact that CAPL relies on the base classes to generate the prototypes for novel categories. Thus, the model may give excessive importance to base classes, not fully leveraging the support samples during the adaptation stage. Furthermore, differences with respect to BAM are also significant, with 7% and 26% improvement on novel classes in the 1-shot and 5-shot settings on PASCAL-5<sup>i</sup>. Note that BAM fuses the output on base classes to the novel prediction and it does not form a holistic classifier over all the classes [14]. Therefore, its prediction on base classes remains intact after learning the novel class, and this leads to high performance on *Base* metric. However, we believe a rigid base prediction hampers the ability of the model to grasp a universal view of the classes, and this comes at the price of worse performance on novel classes, despite the fact that learning them is the main objective of any few-shot learning framework.

The same trend observed in the PASCAL-5<sup>i</sup> benchmark is repeated for COCO-20<sup>i</sup>. More concretely, our method achieves performance gains of around 10% and 17% on *Novel* metric over CAPL, and 3% and 12% compared to BAM under the 1-shot and 5-shot scenarios, respectively.

				PASCAL-5 <sup>i</sup>					
				1-Shot			5-Shot		
Method		Practical setting	Multi-class design	Base	Novel	Mean	Base	Novel	Mean
CANet* [37]	CVPR'19	✓	✗	8.73	2.42	5.58	9.05	1.52	5.29
PANet* [33]	ICCV'19	✓	✗	31.88	11.25	21.57	32.95	15.25	24.1
PFENet* [30]	TPAMI'20	✓	✗	8.32	2.67	5.50	8.83	1.89	5.36
MiB <sup>†</sup> [4]	CVPR'20	✓	✓	63.80	8.86	36.33	68.60	28.93	48.77
SCL* [36]	CVPR'21	✓	✗	8.88	2.44	5.66	9.11	1.83	5.47
RePRI <sup>†</sup> [1]	CVPR'21	✓	✗	20.76	10.50	15.63	34.06	20.98	27.52
CAPL <sup>†</sup> [29]	CVPR'22	✗	✓	64.80	17.46	41.13	65.43	24.43	44.93
BAM <sup>†</sup> [14]	CVPR'22	✗	✗	<b>71.60</b>	27.49	49.55	<b>71.60</b>	28.96	50.28
DIaM	(Ours)	✓	✓	70.89	<b>35.11</b>	<b>53.00</b>	70.85	<b>55.31</b>	<b>63.08</b>
DIaM-UB	(Ours)	✗	✓	71.13	52.61	61.87	71.12	66.12	68.62

				COCO-20 <sup>i</sup>					
				Base	Novel	Mean	Base	Novel	Mean
RePRI <sup>†</sup> [1]	CVPR'21	✓	✗	5.62	4.74	5.18	8.85	8.84	8.85
CAPL <sup>†</sup> [29]	CVPR'22	✗	✓	43.21	7.21	25.21	43.71	11.00	27.36
BAM <sup>†</sup> [14]	CVPR'22	✗	✗	<b>49.84</b>	14.16	32.00	<b>49.85</b>	16.63	33.24
DIaM	(Ours)	✓	✓	48.28	<b>17.22</b>	<b>32.75</b>	48.37	<b>28.73</b>	<b>38.55</b>
DIaM-UB	(Ours)	✗	✓	48.55	29.48	39.02	48.63	40.43	44.53

Table 1. **Quantitative evaluation on PASCAL-5<sup>i</sup> and COCO-20<sup>i</sup> compared to FSS and GFSS methods.** DIaM represents our method and DIaM-UB is an impractical extension of it, explained in Sec. 5.3. All the methods employ ResNet-50 as backbone. Results with a “\*” sign are obtained from [29], whereas results with a “†” sign are reproduced using the publicly available codes.

**Impact of increasing the number of novel classes.** A more challenging scenario involves expanding the set of novel classes. To this end, we compare our approach to CAPL [29] and BAM [14] on the newly defined PASCAL-10<sup>i</sup>, which contains 10 base and 10 novel classes. Table 2 shows that the difference in *Novel* scores compared to these methods is further exacerbated. In particular, DIaM yields improvements on novel classes of nearly 15% and 30% under 1-shot and 5-shot settings, respectively, while also outperforming both methods on base classes.

		1-Shot			5-Shot		
Method		Base	Novel	Mean	Base	Novel	Mean
CAPL [29]	CVPR'22	53.78	15.01	34.40	57.02	20.40	38.71
BAM [14]	CVPR'22	69.02	15.48	42.25	69.18	21.51	45.35
DIaM	(Ours)	<b>70.26</b>	<b>31.29</b>	<b>50.77</b>	<b>70.25</b>	<b>51.89</b>	<b>61.07</b>

Table 2. **Quantitative evaluation on PASCAL-10<sup>i</sup>.** All the methods employ ResNet-50 as backbone.

### 5.3. Ablation studies

We ablate along three axes to better understand each design choice’s contribution. Results are summarized in the form of convergence plots in Fig. 1. Specifically, we find that (a) terms act symbiotically to provide the best performance on both base and novel classes, (b) self-estimation yields higher novel-class performance than the uniform

prior, and finally, (c) as stated in Sec. 4.2.3, introducing the knowledge distillation term and optimizing the base classifier  $\theta_b$ , as opposed to the simple solution of freezing it, allows both faster convergence and better optima. Note that since we form a holistic distribution over all classes, even if we freeze  $\theta_b$ , the probability of base classes will not remain fixed. Those convergence plots demonstrate that models improve rapidly at first and keep improving at a slower pace as the adaptation continues (more details in Appendix B).

Based on Fig. 1, the self-estimation of  $\Pi$  leads to better performance than using the uniform distribution. Therefore, we wonder: *how far could this term take us?* The methods mentioned as DIaM-UB in Tab. 1 demonstrate an extension of our model in which the actual true  $\Pi$  is given. This shows the upper bound of the performance of our model as the estimation of  $\Pi$  gets more accurate. Experiments show that an accurate  $\Pi$  can lead to significant improvements in performance. We emphasize that knowing the exact size of the target objects might be unrealistic and that our goal is just to demonstrate that with a proper mechanism to provide an accurate estimate of the target class proportion, the obtained results can be much improved.

### 5.4. Visual examples

Qualitative results of the proposed method are presented in Fig. 2. More specifically, for a given query, the predictions made by four different models, each containing a sub-

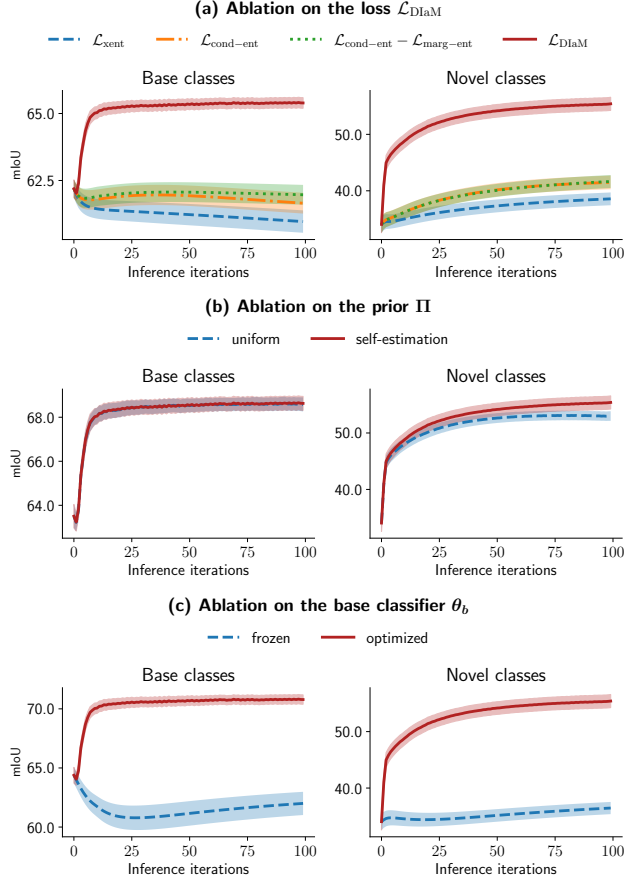


Figure 1. **Ablation studies** on (a) DIaM’s loss, (b) the type of prior  $\Pi$  in Eq. (9), and (c) the optimization of  $\theta_b$ . Results are provided for PASCAL-5<sup>i</sup> under the 5-shot setting.

set of our loss function are shown. This figure reveals an interesting phenomenon: in the absence of the knowledge distillation term of Eq. (12), the model tends to predict some of the base classes as the novel ones. For example, in the third and fourth row, the base classes *horse* and *bicycle* are mistakenly segmented as the novel classes *cow* and *car*. These misclassifications are revised when the knowledge distillation term is in use, which further proves its effectiveness. More visual examples are provided in [Appendix H](#).

## 6. Conclusion

**Summary.** We propose a new generalized few-shot segmentation method, with a standard supervised training scheme and a lightweight inference phase, which can be applied on top of any learned feature extractor and classifier. Our method is based on the InfoMax framework [18] incorporating problem-specific biases, and it also employs knowledge distillation [12] to prevent performance loss on the classes learned during training. Compared to prior

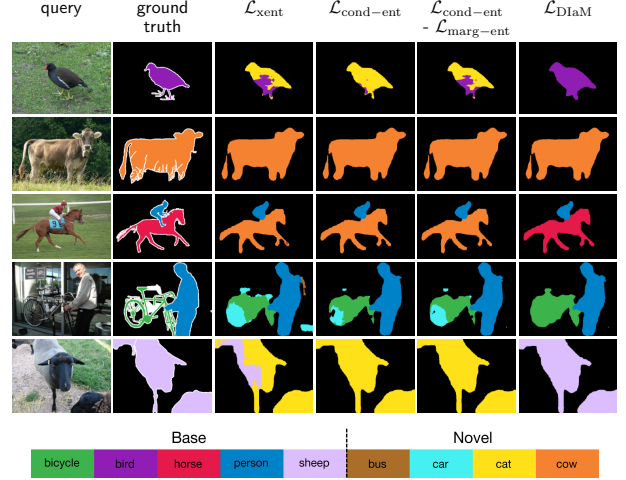


Figure 2. **Qualitative results of different terms of DIaM’s loss function (on PASCAL-5<sup>i</sup> under the 5-shot setting).** A single support set, containing novel classes *bus*, *cat*, *car*, *chair*, and *cow*, is used for predicting every query image. Query images can contain any classes and every one of them is to be recognized.

works, our results show significant improvement in learning novel classes, while keeping the performance on base classes high as well. We eliminated some limiting assumptions of prior methods, such as recognizing one novel class at a time, benefiting from some information about novel classes during training, and having to label base classes in support images. Our proposed knowledge distillation considerably helps retain base knowledge, and we believe imposing such a term is more realistic and practical than explicit supervision for base classes.

**Limitations.** Results of the DIaM-UP experiments show that our marginal entropy term in Eq. (9) can play a significant role and increase the performance considerably. In particular, the ablation study demonstrates that even though our simple choice of estimating the prior proportion  $\Pi$  using the model’s predictions introduces slight improvements, access to a more precise prior has the potential to substantially improve the results. Therefore, we believe that the presented results can be further improved, and encourage future research to explore more powerful mechanisms to provide more accurate proportions of the object of interest.

## Acknowledgements

This work is supported by the National Science and Engineering Research Council of Canada (NSERC), Fonds de recherche du Québec (FRQNT), and Prompt Quebec. We also thank Calcul Quebec and Compute Canada.



## References

- [1] Malik Boudiaf, Hoel Kervadec, Ziko Imtiaz Masud, Pablo Piantanida, Ismail Ben Ayed, and Jose Dolz. Few-shot segmentation without meta-learning: A good transductive inference is all you need? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13979–13988, 2021. [1](#), [2](#), [5](#), [6](#), [7](#)
- [2] Malik Boudiaf, Imtiaz Ziko, Jérôme Rony, José Dolz, Pablo Piantanida, and Ismail Ben Ayed. Information maximization for few-shot learning. *Advances in Neural Information Processing Systems*, 33:2445–2457, 2020. [2](#)
- [3] Francisco M Castro, Manuel J Marín-Jiménez, Nicolás Guil, Cordelia Schmid, and Karteek Alahari. End-to-end incremental learning. In *Proceedings of the European conference on computer vision (ECCV)*, pages 233–248, 2018. [5](#)
- [4] Fabio Cermelli, Massimiliano Mancini, Samuel Rota Buló, Elisa Ricci, and Barbara Caputo. Modeling the background for incremental learning in semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9233–9242, 2020. [5](#), [6](#), [7](#)
- [5] Guneet Singh Dhillon, Pratik Chaudhari, Avinash Ravichandran, and Stefano Soatto. A baseline for few-shot image classification. In *International Conference on Learning Representations*, 2019. [2](#)
- [6] Jiahua Dong, Lixu Wang, Zhen Fang, Gan Sun, Shichao Xu, Xiao Wang, and Qi Zhu. Federated class-incremental learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10164–10173, 2022. [5](#)
- [7] Nanqing Dong and Eric P Xing. Few-shot semantic segmentation with prototype learning. In *BMVC*, volume 3, 2018. [2](#)
- [8] Mark Everingham, SM Eslami, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes challenge: A retrospective. *International journal of computer vision*, 111(1):98–136, 2015. [5](#), [11](#)
- [9] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning*, pages 1126–1135. PMLR, 2017. [2](#)
- [10] Bharath Hariharan, Pablo Arbeláez, Ross Girshick, and Jitendra Malik. Simultaneous detection and segmentation. In *European conference on computer vision*, pages 297–312. Springer, 2014. [5](#), [11](#)
- [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. [6](#)
- [12] Geoffrey Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the knowledge in a neural network. In *NIPS Deep Learning and Representation Learning Workshop*, 2015. [8](#)
- [13] Mohammed Jabi, Marco Pedersoli, Amar Mitiche, and Ismail Ben Ayed. Deep clustering: On the link between discriminative models and k-means. *IEEE transactions on pattern analysis and machine intelligence*, 43(6):1887–1896, 2019. [4](#)
- [14] Chunbo Lang, Gong Cheng, Binfei Tu, and Junwei Han. Learning what not to segment: A new perspective on few-shot segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8057–8067, 2022. [1](#), [2](#), [3](#), [6](#), [7](#), [11](#), [12](#), [13](#)
- [15] Gen Li, Varun Jampani, Laura Sevilla-Lara, Deqing Sun, Jonghyun Kim, and Joongkyu Kim. Adaptive prototype learning and allocation for few-shot segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8334–8343, 2021. [2](#)
- [16] Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2935–2947, 2017. [5](#)
- [17] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. [5](#), [11](#)
- [18] Ralph Linsker. Self-organization in a perceptual network. *Computer*, 21(3):105–117, 1988. [4](#), [8](#)
- [19] Yanbin Liu, Juho Lee, Minseop Park, Saehoon Kim, Eunho Yang, Sung Ju Hwang, and Yi Yang. Learning to propagate labels: Transductive propagation network for few-shot learning. In *International Conference on Learning Representations*, 2019. [2](#)
- [20] Yuanwei Liu, Nian Liu, Xiwen Yao, and Junwei Han. Intermediate prototype mining transformer for few-shot semantic segmentation. In *Neural Information Processing Systems (NeurIPS)*, 2022. [2](#)
- [21] 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. [2](#)
- [22] Zhihe Lu, Sen He, Xiatian Zhu, Li Zhang, Yi-Zhe Song, and Tao Xiang. Simpler is better: Few-shot semantic segmentation with classifier weight transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8741–8750, 2021. [2](#)
- [23] Kate Rakelly, Evan Shelhamer, Trevor Darrell, Alyosha Efros, and Sergey Levine. Conditional networks for few-shot semantic segmentation. In *International Conference on Learning Representations (ICLR) Workshop*, 2018. [2](#)
- [24] Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. In *International Conference on Learning Representations (ICLR)*, 2017. [2](#)
- [25] Amirreza Shaban, Shray Bansal, Liu Zhen, Irfan Essa, and Byron Boots. One-shot learning for semantic segmentation. In *Proceedings of the British Machine Vision Conference (BMVC)*, pages 167.1–167.13. BMVA Press, September 2017. [2](#), [5](#), [11](#)
- [26] Claude Elwood Shannon. A mathematical theory of communication. *The Bell system technical journal*, 27(3):379–423, 1948. [4](#)
- [27] Mennatullah Siam, Boris N Oreshkin, and Martin Jagersand. Amp: Adaptive masked proxies for few-shot segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5249–5258, 2019. [2](#)

- [28] Yanpeng Sun, Qiang Chen, Xiangyu He, Jian Wang, Haocheng Feng, Junyu Han, Errui Ding, Jian Cheng, Zechao Li, and Jingdong Wang. Singular value fine-tuning: Few-shot segmentation requires few-parameters fine-tuning. In *NeurIPS*, 2022. 1, 12
- [29] Zhuotao Tian, Xin Lai, Li Jiang, Shu Liu, Michelle Shu, Hengshuang Zhao, and Jiaya Jia. Generalized few-shot semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11563–11572, 2022. 1, 2, 3, 6, 7, 11, 12
- [30] 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*, 2020. 1, 2, 6, 7
- [31] Olivier Veilleux, Malik Boudiaf, Pablo Piantanida, and Ismail Ben Ayed. Realistic evaluation of transductive few-shot learning. *Advances in Neural Information Processing Systems*, 34:9290–9302, 2021. 5
- [32] Haochen Wang, Xudong Zhang, Yutao Hu, Yandan Yang, Xianbin Cao, and Xiantong Zhen. Few-shot semantic segmentation with democratic attention networks. In *European Conference on Computer Vision*, pages 730–746. Springer, 2020. 2
- [33] Kaixin Wang, Jun Hao Liew, Yingtian Zou, Daquan Zhou, and Jiashi Feng. Panet: Few-shot image semantic segmentation with prototype alignment. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9197–9206, 2019. 1, 2, 6, 7
- [34] Boyu Yang, Chang Liu, Bohao Li, Jianbin Jiao, and Qixiang Ye. Prototype mixture models for few-shot semantic segmentation. In *European Conference on Computer Vision*, pages 763–778. Springer, 2020. 2
- [35] Han-Jia Ye, Hexiang Hu, and De-Chuan Zhan. Learning adaptive classifiers synthesis for generalized few-shot learning. *International Journal of Computer Vision*, 129:1930–1953, 2021. 11
- [36] Bingfeng Zhang, Jimin Xiao, and Terry Qin. Self-guided and cross-guided learning for few-shot segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8312–8321, 2021. 2, 6, 7
- [37] Chi Zhang, Guosheng Lin, Fayao Liu, Rui Yao, and Chunhua Shen. Canet: Class-agnostic segmentation networks with iterative refinement and attentive few-shot learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5217–5226, 2019. 6, 7
- [38] Gengwei Zhang, Guoliang Kang, Yi Yang, and Yunchao Wei. Few-shot segmentation via cycle-consistent transformer. *Advances in Neural Information Processing Systems*, 34:21984–21996, 2021. 2
- [39] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2881–2890, 2017. 6
- [40] Imtiaz Ziko, Jose Dolz, Eric Granger, and Ismail Ben Ayed. Laplacian regularized few-shot learning. In *International conference on machine learning*, pages 11660–11670. PMLR, 2020. 2

## A. Datasets

We evaluated our method on two widely-used few-shot segmentation benchmarks: PASCAL-5<sup>i</sup> [25] and COCO-20<sup>i</sup> [25]. The former is built based on PASCAL VOC 2012 [8] (containing 20 semantic classes) with additional annotations from SDS [10], while the latter is built from MS-COCO [17] (containing 80 semantic classes). In both datasets, the classes are split into 4 disjoint subsets, and the experiments are done in a cross-validation manner. For each fold, the set of novel classes is extracted from one of these subsets while the union of the remaining subsets will be the set of base classes. Furthermore, as discussed in the main paper, we have introduced a new scenario, where the number of novel classes is increased, referred to as PASCAL-10<sup>i</sup>. The semantic classes in each fold of PASCAL-10<sup>i</sup> are detailed in Tab. 3.

## B. Ablation on the number of iterations

In our empirical validation, the proposed loss function,  $\mathcal{L}_{\text{DlaM}}$ , is optimized for a fixed number of iterations ( $n = 100$ ), which was chosen arbitrarily. As demonstrated in Fig. 3 and Tab. 4 the metrics reach high values using only a few iterations. This finding shows that we can speed up the adaptation further by reducing the number of iterations while keeping the performance relatively intact. Also, continuing the adaptation for longer does not hurt the performance and the metrics stay almost the same. Please note that, as stated, the number of iterations in all the experiments in the main manuscript was set to 100 arbitrarily, disregarding the findings of this section.

## C. Detailed results

The evaluation of our approach is performed in a cross-validation manner. In particular, there exist 4 folds for each of PASCAL-5<sup>i</sup> and COCO-20<sup>i</sup> benchmarks, and 2 folds for PASCAL-10<sup>i</sup>. In the main manuscript, the reported results are obtained by averaging over all the folds in these benchmarks. Table 5 shows the performance of our model on each fold individually.

	Novel classes	Base classes
PASCAL-10 <sup>0</sup>	aeroplane, bicycle, bird, boat, bottle, bus, car, cat, chair, cow	diningtable, dog, horse, motorbike, person, potted plant, sheep, sofa, train, tv/monitor
PASCAL-10 <sup>1</sup>	diningtable, dog, horse, motorbike, person, potted plant, sheep, sofa, train, tv/monitor	aeroplane, bicycle, bird, boat, bottle, bus, car, cat, chair, cow

Table 3. **Semantic classes in each fold of PASCAL-10<sup>i</sup>**. In this benchmark, each fold contains 10 novel classes and hence introduces more difficulties in the generalized few-shot segmentation scenario.

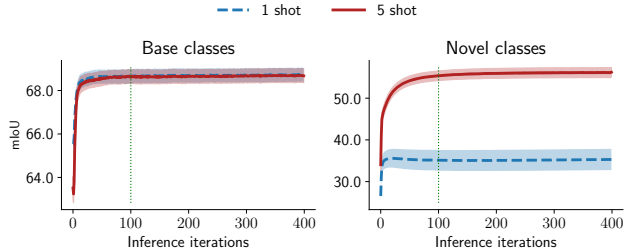


Figure 3. **Performance of the method as the number of iterations in the adaptation phase increases.** Metrics improve rapidly and first and the improvements slow down as the model is further optimized. The dotted green line indicates our choice for the number of iterations in the main manuscript. Results are provided for PASCAL-5<sup>i</sup>.

## D. Harmonic mean

Following [35], we provide in Tab. 6 the harmonic mean score, referred to as *H-Mean*, of CAPL [29], BAM [14], and our method for reference. Using this metric increases the overall performance gap between our method and existing approaches.

## E. Including the background in the base score

CAPL [29] takes into account the background IoU, which is generally higher than the IoU of base classes, when computing the *Base* metric. We, the same as [14], believe that since background does not represent an object of interest, the model’s performance on this class should not be considered. Nevertheless, including the background IoU in the metrics leads to marginal performance differences that are consistent across all methods. In Tab. 7, for GFSS methods, the background IoU is included in the *Base* metric, re-framing it as *Base w/ bg* to avoid confusion.

## F. Practical setting: employing the whole training dataset

As discussed in the main manuscript, CAPL [29] and BAM [14] filter out training images that contain novel classes. This procedure is impractical in real-world scenarios since it needs a training set in which novel classes

# iterations	1-Shot			5-Shot		
	Base	Novel	Mean	Base	Novel	Mean
10	70.15	35.38	52.77	69.92	48.51	59.22
50	70.79	35.33	53.06	70.63	54.15	62.39
100	70.89	35.11	53.00	70.85	55.31	63.08
200	70.87	35.10	52.99	70.81	56.00	63.41
300	70.91	35.18	53.05	70.83	56.19	63.51
400	70.88	35.30	53.09	70.85	56.22	63.54

Table 4. **Precise values of the performance metrics, at selected points on the plots in Fig. 3.** The shaded row indicates our choice for the number of iterations reported in the main manuscript. Results are provided for PASCAL-5<sup>i</sup>.

Benchmark	Fold	1-Shot			5-Shot		
		Base	Novel	Mean	Base	Novel	Mean
PASCAL-5 <sup>i</sup>	0	71.33	29.36	50.35	71.06	53.72	62.39
	1	69.54	46.72	58.13	69.63	63.33	66.48
	2	69.10	27.07	48.09	69.12	54.01	61.57
	3	73.60	37.30	55.45	73.60	50.19	61.90
	mean	70.89	35.11	53.00	70.85	55.31	63.08
COCO-20 <sup>i</sup>	0	49.01	15.89	32.45	48.90	24.86	36.88
	1	46.83	19.50	33.17	47.10	33.94	40.52
	2	48.82	16.93	32.88	49.12	27.15	38.14
	3	48.45	16.57	32.51	48.37	28.95	38.66
	mean	48.28	17.22	32.75	48.37	28.73	38.55
PASCAL-10 <sup>i</sup>	0	68.69	34.40	51.55	68.49	55.94	62.22
	1	71.83	28.17	50.00	72.00	47.84	59.92
	mean	70.26	31.29	50.77	70.25	51.89	61.07

Table 5. **Detailed results for each fold.** For each of the benchmarks, the performance of our method is presented for all the folds.

Method	PASCAL-5 <sup>i</sup>					
	1-Shot			5-Shot		
	Base	Novel	H-Mean	Base	Novel	H-Mean
CAPL [29]	64.80	17.46	27.51	65.43	24.43	35.58
BAM [14]	<b>71.60</b>	27.49	39.73	<b>71.60</b>	28.96	41.24
DiaM	70.89	<b>35.11</b>	<b>46.96</b>	70.85	<b>55.31</b>	<b>62.12</b>
	COCO-20 <sup>i</sup>					
	Base	Novel	H-Mean	Base	Novel	H-Mean
	Base	Novel	H-Mean	Base	Novel	H-Mean
CAPL [29]	43.21	7.21	12.36	43.71	11.00	17.58
BAM [14]	<b>49.84</b>	14.16	22.05	<b>49.85</b>	16.63	24.94
DiaM	48.28	<b>17.22</b>	<b>25.39</b>	48.37	<b>28.73</b>	<b>36.05</b>

Table 6. **Quantitative evaluation on PASCAL-5<sup>i</sup> and COCO-20<sup>i</sup> compared to GFSS methods, using harmonic mean as the overall score.**

are labeled, undermining the goal of few-shot learning, *i.e.*, having only a few labeled examples of the novel classes. Recent empirical evidence [28] has shown that such additional step can lead to performance gain on novel classes. In Tab. 8, we have changed the training procedure of CAPL and BAM and avoided removing images containing novel classes from training. More specifically, the potential ob-

Method	1-Shot			5-Shot		
	Base w/ bg	Novel	Mean	Base w/ bg	Novel	Mean
CAPL [29]	66.37	17.46	41.92	66.95	24.43	45.69
BAM [14]	72.00	27.49	49.75	<b>72.36</b>	28.96	50.66
DiaM	<b>72.04</b>	<b>35.11</b>	<b>53.58</b>	72.12	<b>55.31</b>	<b>63.72</b>

Table 7. **Quantitative evaluation on PASCAL-5<sup>i</sup> compared to GFSS methods, including the background performance in the metrics.**

jects from novel classes are labeled as *background* during training.

Method	1-Shot			5-Shot		
	Base	Novel	Mean	Base	Novel	Mean
CAPL [29]	71.59	12.69	42.14	<b>71.71</b>	19.58	45.65
BAM [14]	<b>71.61</b>	19.35	45.48	71.66	26.33	49.00
DiaM	70.89	<b>35.11</b>	<b>53.00</b>	70.85	<b>55.31</b>	<b>63.08</b>

Table 8. **Quantitative evaluation on PASCAL-5<sup>i</sup> compared to GFSS methods, in the experimental setting in which the whole training dataset is employed.** In this setting, images containing novel classes are not removed from the training process. Confirming the findings in [28], this procedure enhances the performance on novel classes. It is also worth noting that although *Novel* score has been decreased for CAPL, its *Base* score has been considerably increased.

## G. Adaptation of BAM to multi-class GFSS

The results reported in [14] for the GFSS task are based on an evaluation protocol in which only one novel class can be recognized in a query image. Indeed, the *meta-learner* from this method can only provide binary (*i.e.*, *background vs foreground*) predictions and is not practical in the setting where multiple novel classes are to be predicted at the same time. To be able to incorporate BAM in our empirical validation, we had to adapt it so that it can predict multiple novel classes simultaneously. These modifications are detailed in what follows. First, instead of selecting  $K$  support samples of a novel class  $c$  and asking the model to segment class  $c$  in a query image, we form  $|\mathcal{C}^n|$  different support sets,  $\mathbb{S}_i$ , one for each  $i \in \mathcal{C}^n$ . Recall that  $\mathcal{C}^n$  is the set of novel classes and that  $\mathbb{S}_i$  contains  $K$  samples labeled for each novel class  $i$ . Second, we run BAM  $|\mathcal{C}^n|$  times and, in each inference, we give the same query image alongside  $\mathbb{S}_i$ , resulting in a foreground probability map for each class  $i$ , called  $\mathbf{m}_i$ . Then, we need to create a single mask containing all the novel class predictions to further use it in BAM’s fusion mechanism. To do this, we create an aggregated novel map,  $\mathbf{a}$ , which is formed based on the resulted  $|\mathcal{C}^n|$  maps,



in such a way that for each pixel  $j$ :

$$\mathbf{a}(j) = \operatorname{argmax}_{i \in \mathcal{C}^n} \mathbf{m}_i(j). \quad (14)$$

We also form  $\mathbf{p}_\mathbf{a}$  to preserve the probability of the selected indices, which will later be compared to the predefined threshold  $\tau$  introduced in [14]:

$$\mathbf{p}_\mathbf{a}(j) = \max_{i \in \mathcal{C}^n} \mathbf{m}_i(j). \quad (15)$$

Then,  $\mathbf{a}$  and  $\mathbf{p}_\mathbf{a}$  alongside the base map predicted by BAM’s *base-learner* for the query,  $\hat{\mathbf{m}}_b$ , are used to perform the fusion procedure following [14]. More specifically, the final prediction is formulated as

$$\hat{\mathbf{m}}_g(j) = \begin{cases} \mathbf{a}(j) & \mathbf{p}_\mathbf{a}(j) > \tau, \\ \hat{\mathbf{m}}_b(j) & \mathbf{p}_\mathbf{a}(j) \leq \tau \text{ and } \hat{\mathbf{m}}_b(j) \neq 0, \\ 0 & \text{otherwise.} \end{cases} \quad (16)$$

This change definitely slows the inference by an order of magnitude, but this is inevitable because of the nature of meta-learning few-shot segmentation, which needs to be accommodated to produce multi-class semantic maps, as they are tailored to binary maps.

## H. Visual examples

In the main manuscript, we presented qualitative results on PASCAL-5<sup>i</sup> using different versions of our loss function. We observed that in the absence of the knowledge distillation term, the model misclassifies some of the previously learned base classes as novel ones. Figure 4 shows similar results on COCO-20<sup>i</sup>, where the same trend is observed. For instance, in the first two rows, base classes *cell phone* and *keyboard* are mistakenly classified as the novel class *remote*. Note that this problem is fixed when the knowledge distillation term is added to the loss function.

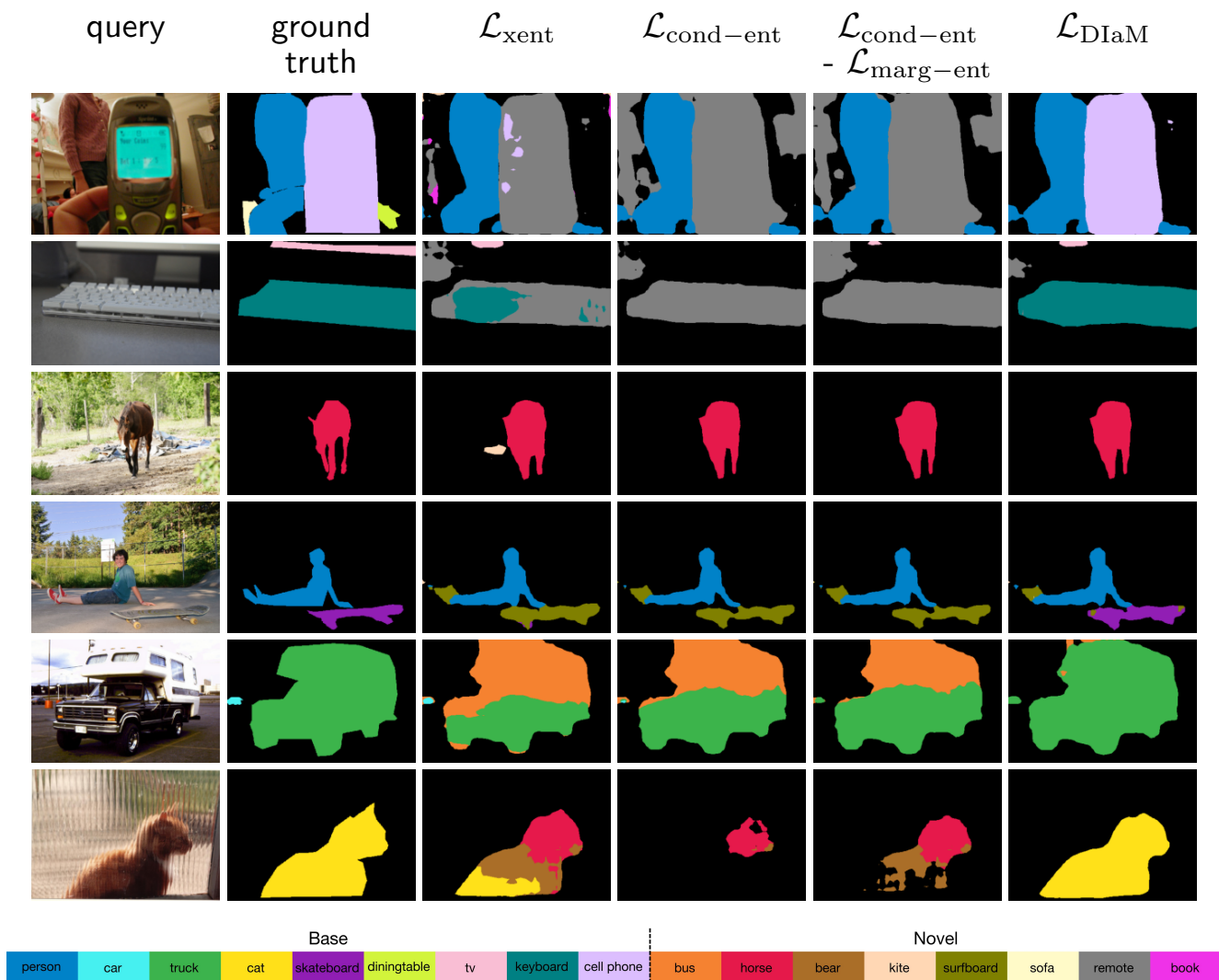


Figure 4. **Qualitative results of different terms of DIaM’s loss function on COCO-20<sup>i</sup>.** A single support set, containing the following novel classes is used for predicting every query image: *bicycle, bus, traffic light, bench, horse, bear, umbrella, frisbee, kite, surfboard, cup, bowl, orange, pizza, sofa, toilet, remote, oven, book, teddy bear*. Query images can contain any classes and every one of them is to be recognized. From the left, the first two columns show the query image and the ground truth, and the following columns display predictions of models using different loss functions. Results are on COCO-20<sup>i</sup> under the 5-shot setting.