

Efficient Masked Attention Transformer for Few-Shot Classification and Segmentation

Dustin Carrión-Ojeda^{*1,2} , Stefan Roth^{1,2} , and Simone Schaub-Meyer^{1,2} 

¹ Department of Computer Science, Technical University of Darmstadt, Germany

² Hessian Center for AI (hessian.AI), Germany

{dustin.carrión, stefan.roth, simone.schaub}@visinf.tu-darmstadt.de

<https://visinf.github.io/emat>

Abstract. Few-shot classification and segmentation (FS-CS) focuses on jointly performing multi-label classification and multi-class segmentation using few annotated examples. Although the current state of the art (SOTA) achieves high accuracy in both tasks, it struggles with small objects. To overcome this, we propose the **Efficient Masked Attention Transformer** (EMAT), which improves classification and segmentation accuracy, especially for small objects. EMAT introduces three modifications: a novel memory-efficient masked attention mechanism, a learnable downscaling strategy, and parameter-efficiency enhancements. EMAT outperforms all FS-CS methods on the PASCAL-5ⁱ and COCO-20ⁱ datasets, using at least four times fewer trainable parameters. Moreover, as the current FS-CS evaluation setting discards available annotations, despite their costly collection, we introduce two novel evaluation settings that consider these annotations to better reflect practical scenarios.

Keywords: Few-shot learning · Efficiency · Segmentation · Classification.

1 Introduction

Recently, data-intensive methods have been introduced for various deep learning applications [5, 8, 23, 25, 32, 34, 41]. These methods rely on large training datasets, making them impractical in fields where collecting extensive datasets is challenging or costly [13, 14, 65]. Consequently, few-shot learning (FSL) methods have gained significant attention for their ability to learn from just a few examples and quickly adapt to new classes [1, 44, 52, 56]. In computer vision, FSL has been mostly applied to image classification (FS-C) [3, 18, 40, 43] and segmentation (FS-S) [11, 30, 55, 62, 63].

FS-C and FS-S often co-occur in real-world applications, *e.g.*, in agriculture, where crops must be segmented and classified by type or health status. Hence, recent works [19, 20] integrate multi-label classification and multi-class segmentation into a single few-shot classification and segmentation (FS-CS) task. While

* Corresponding author.



Fig. 1. Qualitative comparison of small objects (*i.e.*, objects that occupy less than 15% of the image) between the current SOTA FS-CS method (CST) [20] and our proposed EMAT. CST* uses the same backbone as EMAT (*i.e.*, DINOv2 [32]). By processing high-resolution correlation tokens, EMAT preserves finer details, yielding more accurate segmentation masks.

FS-CS addresses some limitations of FS-C (*e.g.*, assuming the query image contains only one class) and FS-S (*e.g.*, assuming the target class is always present in the query image), it also increases the task difficulty by simultaneously tackling classification and segmentation. Moreover, some applications, *e.g.*, medical imaging, rely on precise small-object analysis [13, 16, 65]. Thus, achieving high accuracy on small objects is a desired property for FS-CS methods. Yet, as shown in Fig. 1, the current state-of-the-art (SOTA) FS-CS method [20] struggles with small objects, a limitation we address in this work.

To better align the evaluation of FS-CS models with practical scenarios, FS-CS uses the N -way K -shot configuration, where the model learns N classes from $N \times K$ examples (K per class). However, the current evaluation setting [19, 20] discards available annotations, which is not ideal given the cost of data annotation. To address this, we introduce two new evaluation settings.

Contributions. (1) Building on the current SOTA FS-CS method [20], we propose an efficient masked attention transformer (EMAT), which enhances classification and segmentation accuracy, particularly for small objects, while using approximately four times fewer trainable parameters. (2) Our EMAT outperforms all FS-CS methods on the PASCAL-5ⁱ and COCO-20ⁱ datasets, supports the N -way K -shot configuration, and can generate empty segmentation masks when no target objects are present. (3) Finally, we introduce two new FS-CS evaluation settings that better utilize available annotations during inference.

2 Related Work

Few-shot Classification (FS-C) methods can be categorized into three groups based on what the model learns. *Representation-based* approaches learn class-agnostic, discriminative embeddings [3, 17, 21, 39, 46, 48, 60]. *Optimization-based* approaches learn the optimal set of weights that allow the model to adapt to new classes in just a few optimization steps [4, 15, 35, 40]. *Transfer-based* approaches adapt large pre-trained [6, 10, 24, 28, 43] or foundation models [18, 37, 64]. A major limitation of most FS-C methods is the assumption of a single label per image [2, 38], limiting them in multi-label settings.

Few-shot Segmentation (FS-S) methods can also be categorized into three groups: *prototype matching*, which aligns support embeddings with query features [11, 26, 45, 50, 51, 58]; *dense correlation*, which constructs support–query correlation tensors [7, 29, 30, 33, 53, 54]; and *model-adaptation*, which fine-tunes large pre-trained models [27, 49, 57, 61, 62]. Despite the advancements in FS-S, most methods have two main limitations: (1) they target only the 1-way K -shot configuration and (2) they assume the query image contains the target class, preventing the models from predicting empty segmentation masks. Only a few recent works [42, 59] address the more general N -way K -shot configuration.

Few-shot Classification and Segmentation (FS-CS) focuses on jointly predicting the multi-label classification vector and multi-class segmentation mask without assuming support classes are present in the query image [19]. The current SOTA FS-CS method, the classification-segmentation transformer (CST) [20], uses a memory-intensive masked-attention mechanism that requires significant downsampling of the correlation features, reducing its accuracy on small objects. In this work, we enhance CST by proposing an efficient masked-attention formulation and adding further refinements, resulting in a more memory- and parameter-efficient method with improved accuracy, especially for small objects.

3 Problem Definition

This work focuses on the few-shot classification and segmentation (FS-CS) task [19], formulated as an N -way K -shot learning problem [48]. We assume two disjoint class sets: $\mathcal{C}_{\text{train}}$ for training and $\mathcal{C}_{\text{test}}$ for testing. Accordingly, training tasks are sampled from $\mathcal{C}_{\text{train}}$, and testing tasks from $\mathcal{C}_{\text{test}}$. Each task consists of a support set \mathcal{S} and a query image \mathbf{I}_q , where \mathcal{S} contains N classes \mathcal{C}_s ($\mathcal{C}_s \subseteq \mathcal{C}_{\text{train}}$ or $\mathcal{C}_s \subseteq \mathcal{C}_{\text{test}}$), each represented by K examples:

$$\mathcal{S} = \left\{ \left\{ (\mathbf{I}_j^i, \mathbf{M}_j^i, y_j^i) \mid y_j^i \in \mathcal{C}_s \right\}_j^K \right\}_i^N, \quad (1)$$

where \mathbf{I}_j^i , \mathbf{M}_j^i , and y_j^i denote the support image, segmentation mask, and class label for the j^{th} example of the i^{th} class. Although $y_j^i = i \ \forall j$ in Eq. (1), we use this notation for compatibility with multi-label settings where \mathbf{y}_j^i can vary.

The goal of FS-CS is to learn from \mathcal{S} such that, given \mathbf{I}_q , the model can (i) identify which support classes are present (multi-label classification), and (ii) segment those classes (multi-class segmentation). Moreover, FS-CS allows \mathbf{I}_q to contain a subset of the support classes. Thus, when $N > 1$, \mathbf{I}_q can contain: (1) none of the support classes ($\mathcal{C}_q = \emptyset$), (2) a subset of them ($\mathcal{C}_q \subset \mathcal{C}_s$), or (3) all support classes ($\mathcal{C}_q = \mathcal{C}_s$). Note that case (1) is important in real-world applications where query images may not contain relevant classes, requiring models to predict empty segmentation masks when necessary.

The drawback of the current FS-CS setting is that each support image \mathbf{I}_j^i is assumed to contain only one annotated class (y_j^i). If \mathbf{I}_j^i includes multiple

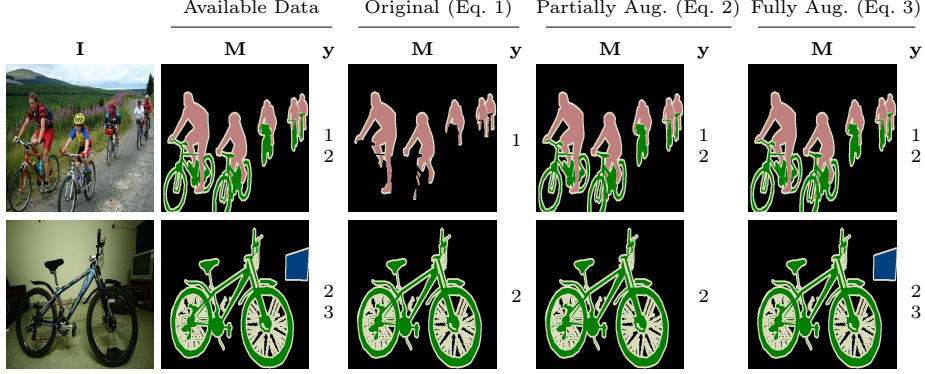


Fig. 2. Example of a 2-way 1-shot (base configuration) support set across different few-shot evaluation settings. \mathbf{I} , \mathbf{M} , and \mathbf{y} represent the images, segmentation masks, and labels, respectively.

support classes ($\mathbf{y}_j^i \subseteq \mathcal{C}_s$), its label vector and segmentation mask need to be adjusted before constructing \mathcal{S} . This adjustment discards available annotations, as illustrated in Fig. 2, where \mathbf{I}_1^1 contains both support classes (person, bike). However, since it is an example of the 1st class (person), the annotations of the 2nd class (bike) are removed in the original setting.

3.1 Proposed Evaluation Settings

To better utilize available annotations and reflect more realistic evaluation scenarios, we introduce two novel FS-CS evaluation settings.

Partially Augmented Setting. This setting keeps all annotations from the support classes:

$$\mathcal{S} = \left\{ \left\{ (\mathbf{I}_j^i, \mathbf{M}_j^i, \mathbf{y}_j^i) \mid \mathbf{y}_j^i \subseteq \mathcal{C}_s \right\}_j^K \right\}_i^N. \quad (2)$$

Note that when $N=1$, this setting is equivalent to Eq. (1). Fig. 2 shows an example for this setting, where \mathbf{M}_1^1 and \mathbf{y}_1^1 keep the annotations of the 2nd class (bike), even though the image is selected as an example of the 1st class (person).

Fully Augmented Setting. This setting keeps all available annotations for each support image, regardless of whether the corresponding classes are part of the support classes:

$$\mathcal{S} = \left\{ \left\{ (\mathbf{I}_j^i, \mathbf{M}_j^i, \mathbf{y}_j^i) \mid \mathbf{y}_j^i \subseteq \mathcal{C}_{\text{train}} \cup \mathcal{C}_{\text{test}} \right\}_j^K \right\}_i^N. \quad (3)$$

For example, in Fig. 2, \mathbf{M}_1^1 and \mathbf{y}_1^1 include annotations for both support classes (person, bike), while \mathbf{M}_1^2 and \mathbf{y}_1^2 are augmented with annotations of a non-support class (TV), which can belong to either $\mathcal{C}_{\text{train}}$ or $\mathcal{C}_{\text{test}}$. In this setting,

the model is expected to classify and segment all support and augmented classes present in \mathbf{I}_q . Additionally, this setting aligns closely with the generalized few-shot setting (GFSL) [44], which also evaluates on base classes. However, unlike standard GFSL, which evaluates on all base classes seen during training, our setting restricts evaluation to only those classes present in the support set.

4 Efficient Masked Attention Transformer

Fig. 3 illustrates the pipeline used by our proposed efficient masked attention transformer (EMAT), which builds upon the classification-segmentation transformer (CST) [20]. Both methods share the same feature extraction process: support and query images $\mathbf{I}_s^i, \mathbf{I}_q \in \mathbb{R}^{H \times W \times 3}$ are processed by a frozen, pre-trained ViT [12] with patch size p , producing support and query image tokens $\mathbf{T}_{s_i}, \mathbf{T}_{q_i} \in \mathbb{R}^{h \times w \times d}$, and a support class token $\mathbf{T}_{s_c} \in \mathbb{R}^{1 \times d}$, where $h = H/p$, $w = W/p$, and d is the token dimension of a single ViT head. The support tokens \mathbf{T}_{s_i} are downsampled via bilinear interpolation and reshaped to $\mathbf{T}_{s_i}^f \in \mathbb{R}^{(h' \cdot w') \times d}$. Similarly, query image tokens \mathbf{T}_{q_i} are reshaped to $\mathbf{T}_{q_i}^f \in \mathbb{R}^{(h \cdot w) \times d}$. Next, $\mathbf{T}_{s_i}^f$ and \mathbf{T}_{s_c} are concatenated to form \mathbf{T}_s^c . Finally, cosine similarity between \mathbf{T}_s^c and $\mathbf{T}_{q_i}^f$ is computed across all ViT layers l and attention heads g , resulting in the correlation tokens $\mathbf{C} \in \mathbb{R}^{t_s \times t_q \times (l \cdot g)}$, where $t_s = h' \cdot w' + 1$ and $t_q = h \cdot w$.

EMAT differs from CST in the two-layer transformer (purple blocks in Fig. 3) that processes the correlation tokens \mathbf{C} and feeds task-specific heads for multi-label classification and multi-class segmentation. EMAT enhances this transformer with three key improvements: (1) a novel memory-efficient masked attention formulation (see Sec. 4.1) that allows using higher-resolution correlation tokens, (2) a learnable downscaling strategy (see Sec. 4.2) that avoids reliance on large pooling kernels, and (3) additional modifications for improved parameter efficiency (see Sec. 4.3), which can help to reduce overfitting a small support set.

Following CST, EMAT is trained using the 1-way 1-shot configuration. Since EMAT uses task-specific heads, it is trained with two losses:

$$\mathcal{L}_{\text{clf}} = -y \log \hat{y}, \quad (4)$$

$$\mathcal{L}_{\text{seg}} = -\frac{1}{HW} \sum_{i=1}^H \sum_{j=1}^W \mathbf{M}_{ij} \log \widehat{\mathbf{M}}_{ij}, \quad (5)$$

where $y \in \{0, 1\}$ and $\mathbf{M}_{ij} \in \{0, 1\}$ are the ground-truth classification and segmentation labels, and \hat{y} , $\widehat{\mathbf{M}}_{ij}$ are the corresponding predictions. The final loss function jointly optimizes both losses using a balancing hyperparameter λ :

$$\mathcal{L} = \lambda \mathcal{L}_{\text{clf}} + \mathcal{L}_{\text{seg}}. \quad (6)$$

Inference on N -way K -shot task is performed as in CST [20], by treating each class as an independent 1-way K -shot task: class-wise logits and segmentation masks are averaged over the K examples, producing N predictions. Logits above a threshold $\delta = 0.5$ form the multi-label vector, and each pixel $\widehat{\mathbf{M}}_{ij}$ is assigned

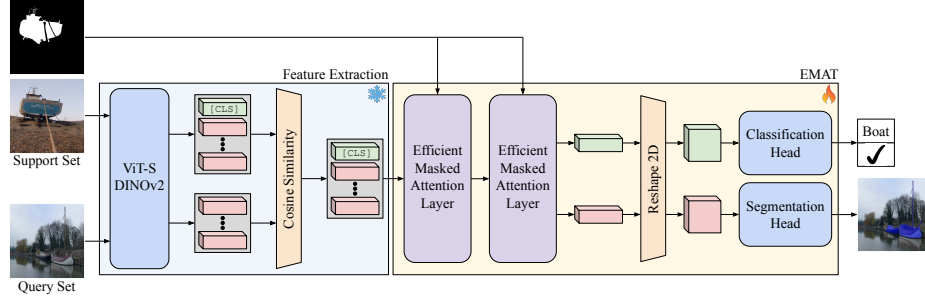


Fig. 3. FS-CS pipeline used by our EMAT. A frozen, pre-trained ViT [12] extracts image and class tokens from support and query images, which are correlated via cosine similarity. The resulting correlation tokens are processed by a two-layer transformer equipped with our masked attention mechanism, learnable downscaling, and parameter-efficient design (see Secs. 4.1 to 4.3). Task-specific heads then predict the multi-label classification vector and multi-class segmentation mask.

to the class with the highest score, or to background if all scores fall below δ , thereby allowing empty masks.

4.1 Memory-Efficient Masked Attention

The SOTA FS-CS method, CST [20], uses a masked attention mechanism based on self-attention [12, 47] to process the correlation tokens $\mathbf{C} \in \mathbb{R}^{t_s \times t_q \times (l \cdot g)}$. Due to the high memory cost of self-attention, CST significantly downsamples the support tokens \mathbf{T}_{s_i} when computing \mathbf{C} , sacrificing fine-grained spatial details (see Fig. 5). To address this, we propose a novel memory-efficient masked attention formulation that allows EMAT to use high-resolution correlation tokens.

Given \mathbf{C} , let $\mathbf{Q}^d \in \mathbb{R}^{t_d \times t_q \times e}$, $\mathbf{K}, \mathbf{V} \in \mathbb{R}^{t_s \times t_q \times e}$ denote the query, key, and value matrices, where e is the embedding size. The dimensions of \mathbf{Q}^d differ from those of \mathbf{K} and \mathbf{V} because the query matrix is progressively downsampled (see Sec. 4.2). Additionally, as shown in Fig. 3, the segmentation mask $\mathbf{M} \in \mathbb{R}^{H \times W}$ enters directly into the attention mechanism without being processed by the feature extractor. Thus, it is resized, flattened, and then append a trailing “1” to obtain $\mathbf{M}^f \in \mathbb{R}^{h' \cdot w' + 1}$. This appended “1” ensures that the support class token is never masked in Eqs. (7) and (8).

For one attention head, CST computes:

$$\mathbf{O}_{ijk} = \sum_p [\text{softmax}((\mathbf{Q}_{ijk}^d \cdot \mathbf{K}_{:jk}) \odot \mathbf{M}_{:}^f)]_p \odot \mathbf{V}_{pjk}, \quad (7)$$

where $i \in \{1, \dots, t_d\}$, $j \in \{1, \dots, t_q\}$, $k \in \{1, \dots, e\}$, and $p \in \{1, \dots, t_s\}$. Here $t_d = h'' \cdot w'' + 1$, $t_q = h \cdot w$, and $t_s = h' \cdot w' + 1$, the “+1” accounts for the support class token, and $(\cdot)'$ denotes a downsampled value. The operator \odot represents element-wise multiplication.

In CST, the support dimension t_s is 145 ($h' = w' = 12$) in the first attention layer and 10 ($h' = w' = 3$) in the second. Consequently, in the second layer

most values of \mathbf{M}^f are zero (see Fig. 5). As a result, Eq. (7) zeros out most attention values, but they are still processed in all intermediate computations, leading to a memory-inefficient formulation. To overcome this, we introduce a memory-efficient reformulation that excludes the masked-out entries:

$$\mathbf{O}_{ijk} = \sum_{p^\circ} [\text{softmax}(\mathbf{Q}_{ijk}^d \cdot (\mathbf{K}_{:jk} \odot \mathbf{M}_{:}^f))]_{p^\circ} \odot (\mathbf{V}_{:jk} \odot \mathbf{M}_{:}^f)_{p^\circ}, \quad (8)$$

where \odot is our element-wise masking operator:

$$(\mathbf{Z}_{pjk} \odot \mathbf{M}_p^f) = \begin{cases} \mathbf{Z}_{pjk} & \text{if } \mathbf{M}_p^f = 1, \\ \emptyset & \text{otherwise,} \end{cases} \quad \forall p \in \{1, \dots, t_s\}, \quad (9)$$

with $\mathbf{Z} \in \mathbb{R}^{t_s \times t_q \times e}$ and \emptyset indicating that the corresponding entry is excluded. This exclusion of elements results in the reduced set of indices $p^\circ \subseteq p$ used in Eq. (8), where $p^\circ = p$ only if \mathbf{M}^f contains no zeros. Excluding masked-out tokens reduces memory usage allowing EMAT to increase the support dimension t_s to 401 ($h' = w' = 20$) in the first attention layer (≈ 2.7 times more than CST) and 101 ($h' = w' = 10$) in the second (≈ 11 times more than CST). Note that the index set p° varies across images in a batch; thus, Eq. (8) is computed sequentially for each image. However, batch processing is still used both before and after this step because the input and output tensors (\mathbf{C} and \mathbf{O}) have the same dimensions for every image. By computing attention only over unmasked entries, EMAT still achieves a runtime comparable to CST.

4.2 Learnable Downscaling

As mentioned in Sec. 4.1, the query matrix $\mathbf{Q} \in \mathbb{R}^{t_s \times t_q \times e}$ is progressively down-scaled. This downscaling occurs before computing the masked-attention in each layer and it keeps t_q and e fixed, while shrinking the support spatial dimensions h' and w' , which reduces the support dimension ($t_s = h' \cdot w' + 1$). Unlike CST, which uses only average pooling, EMAT introduces a lightweight, learnable strategy that combines small convolutions with pooling. This hybrid design removes the need for the large pooling kernels that would otherwise be required to handle the higher-resolution correlation tokens used by EMAT.

In the first attention layer, EMAT splits \mathbf{Q} into support image tokens \mathbf{Q}_i and a single class token \mathbf{Q}_c . After reshaping \mathbf{Q}_i to its $h' \times w'$ support spatial layout, a 3D convolution followed by another reshape produces $\mathbf{Q}_i^r \in \mathbb{R}^{(h'' \cdot w'') \times t_q \times e}$. Simultaneously, a 2D convolution transforms \mathbf{Q}_c into $\mathbf{Q}_c^r \in \mathbb{R}^{1 \times t_q \times e}$. These outputs are concatenated to form the downscaled query matrix $\mathbf{Q}^d \in \mathbb{R}^{t_d \times t_q \times e}$ used in Eq. (8), where $t_d = h'' \cdot w'' + 1$.

The second attention layer repeats the process, but before the concatenation \mathbf{Q}_i^r is collapsed to a single spatial token by a 3D average pool, producing $\mathbf{Q}_i^p \in \mathbb{R}^{1 \times t_q \times e}$. Concatenating \mathbf{Q}_i^p with \mathbf{Q}_c^r gives the downscaled query matrix $\mathbf{Q}^d \in \mathbb{R}^{2 \times t_q \times e}$ used during the attention computation.

4.3 Modifications for Parameter Efficiency

Few-shot models with too many parameters risk overfitting the small support set, thereby reducing their ability to adapt to new classes. Therefore, EMAT reduces the number of channels across all operations: its two attention layers use 64 and 32 channels, versus 32 and 128 in CST, and its two task-specific heads use 32 and 16 channels, versus 128 and 64 in CST. These channel reductions significantly decrease the number of trainable parameters in EMAT.

5 Experiments

Datasets. We evaluated our EMAT on the widely used PASCAL-5ⁱ [36] and COCO-20ⁱ [31] datasets. Although they were designed for few-shot segmentation, both can also be used for few-shot classification and segmentation [19]. PASCAL-5ⁱ comprises 20 classes and COCO-20ⁱ 80 classes, each partitioned into four non-overlapping folds.

Implementation Details. EMAT uses a frozen ViT-S encoder [12] pre-trained with DINOv2 [32]. The two-layer transformer uses our memory-efficient masked attention with 8 heads. We train for 80 epochs with a batch size of 9 using the Adam optimizer [22] with learning rate 10^{-3} . Following [20], we use 1-way 1-shot tasks with the original setting (see Eq. 1) and set the loss weight λ in Eq. (6) to 0.1. Moreover, we re-train CST [20] with the same DINOv2 backbone used by EMAT and denote it as CST*. All training was conducted on three NVIDIA RTX A6000 GPUs, with evaluation performed on a single GPU.

5.1 Comparison to SOTA FS-CS

To evaluate the effectiveness of our EMAT, we compare it with CST [20] and other state-of-the-art few-shot classification and segmentation (FS-CS) methods. Tab. 1 shows the mean classification accuracy (Acc.) and mean Intersection over Union (mIoU) over the four folds of PASCAL-5ⁱ [36] and COCO-20ⁱ [31], for 2-way 1-shot tasks across all evaluation settings (see Sec. 3). Although DINOv2 pre-training [32] already significantly improves CST* over its original version, EMAT consistently outperforms all methods across all settings. These results validate the benefit of processing higher-resolution correlation tokens enabled by our memory-efficient masked attention (see Sec. 4.1). Moreover, EMAT requires at least four times fewer parameters than CST, making it the most parameter-efficient method among SOTA FS-CS models. The supplementary material provides per-fold results for Tab. 1 and additional results on $\{1, \dots, 5\}$ -way 5-shot tasks, demonstrating the scalability of our method.

The results in Tab. 1 also show that our partially augmented setting slightly improves accuracy and mIoU for most methods, confirming the benefit of better exploiting the available annotations. However, the improvement is marginal, likely because only 242 and 106 out of 4000 tasks are augmented for PASCAL-5ⁱ

Table 1. Comparison of FS-CS methods on PASCAL-5ⁱ and COCO-20ⁱ across all evaluation settings: original, partially augmented, and fully augmented, using 2-way 1-shot tasks (base configuration). CST* and EMAT were trained and evaluated, while other methods were only evaluated using the checkpoints from [19]. CST* uses the same backbone as EMAT (*i.e.*, DINOv2 [32]). All values, except the number of trainable parameters (in millions), are percentages (higher is better). **Highlight** indicates our proposed method. **Bold** and underlined values indicate the best and second best results.

Dataset	Method	Train. Params.	Original		Partially Augmented		Fully Augmented	
			Acc.	mIoU	Acc.	mIoU	Acc.	mIoU
PASCAL-5 ⁱ	PANet [50]	23.51	56.53	37.20	56.93	37.49	55.75	37.25
	PFENet [45]	31.96	39.35	35.57	39.48	35.61	36.88	35.08
	HSNet [29]	2.57	67.27	44.85	67.75	44.72	65.92	44.40
	ASNet [19]	1.32	68.30	47.87	68.62	47.78	66.40	47.58
	CST [20]	<u>0.37</u>	70.37	53.78	70.60	53.81	68.45	53.76
	CST*	<u>0.37</u>	<u>80.58</u>	<u>63.28</u>	<u>80.60</u>	<u>63.23</u>	<u>78.57</u>	<u>63.08</u>
	EMAT	0.09	82.70	63.38	82.92	63.32	81.23	63.24
COCO-20 ⁱ	PANet [50]	23.51	51.30	23.64	51.32	23.78	45.07	23.17
	PFENet [45]	31.96	36.45	23.37	36.50	23.39	29.33	21.61
	HSNet [29]	2.57	62.43	30.58	62.40	30.66	55.15	29.44
	ASNet [19]	1.32	63.05	31.62	63.03	31.64	55.47	30.47
	CST [20]	<u>0.37</u>	64.02	36.23	64.10	36.20	56.30	35.60
	CST*	<u>0.37</u>	78.70	51.47	78.87	51.53	71.18	50.76
	EMAT	0.09	80.07	52.81	80.25	52.82	73.00	51.99

and COCO-20ⁱ, respectively. In contrast, our fully augmented setting lowers accuracy and mIoU for every method, although less significantly for mIoU. As discussed in Sec. 3.1, this setting augments not only the support examples but also includes every class present in the support images, making the tasks harder. This setting augments 1243 and 1515 out of the 4000 tasks for PASCAL-5ⁱ and COCO-20ⁱ, respectively. The augmentations in both of our proposed settings highlight that the original evaluation setting fails to use available annotations.

5.2 Qualitative Results

As explained in Sec. 3, when handling N -way K -shot tasks with $N > 1$, the query image can contain (1) none, (2) some, or (3) all of the support classes. The top part of Fig. 4 shows that EMAT produces more accurate segmentation masks than CST* in these three scenarios, confirming that EMAT can predict empty masks and masks with one or multiple classes. The bottom part of Fig. 4 illustrates the same task across all few-shot evaluation settings. In the original setting, both models segment the 1st class (orange) correctly but make errors on the 2nd class (glass), mistakenly segmenting visually similar objects (salt shaker). In the partially augmented setting, the additional annotations for the 2nd class degrade CST, while EMAT maintains a precise segmentation, though

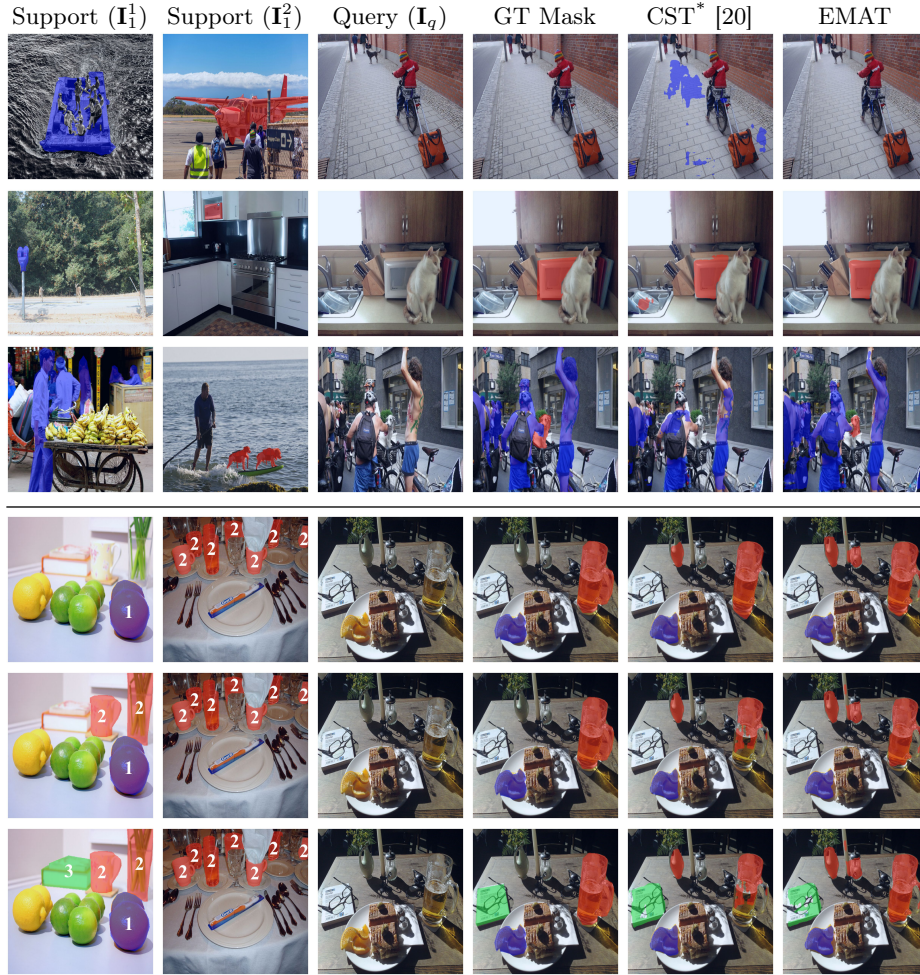


Fig. 4. Qualitative comparison of CST^* vs. our EMAT on COCO-20ⁱ [31] using 2-way 1-shot tasks. **(Top)** Row 1: Query w/o support classes, Row 2: Query w/ subset of support classes, Row 3: Query w/ all support classes. **(Bottom)** Row 1: Original setting, Row 2: Partially augmented setting, Row 3: Fully augmented setting.

both still incorrectly segment non-target objects. In the fully augmented setting, both methods correctly segment the additional class (book).

To illustrate the effect of higher-resolution tokens, Fig. 5 compares the segmentation masks used in the masked-attention layers of CST and EMAT. Thanks to our memory-efficient formulation (see Sec. 4.1), EMAT preserves more details across layers. The difference is most visible in the second layer, where the mask used by CST barely contains any information since it has a resolution of 3×3 , whereas EMAT, with a 10×10 resolution, retains meaningful details and structure. For instance, the person in the bottom-right corner of the last row of Fig. 5 remains visible in both layers of EMAT but vanishes in the second layer of CST.

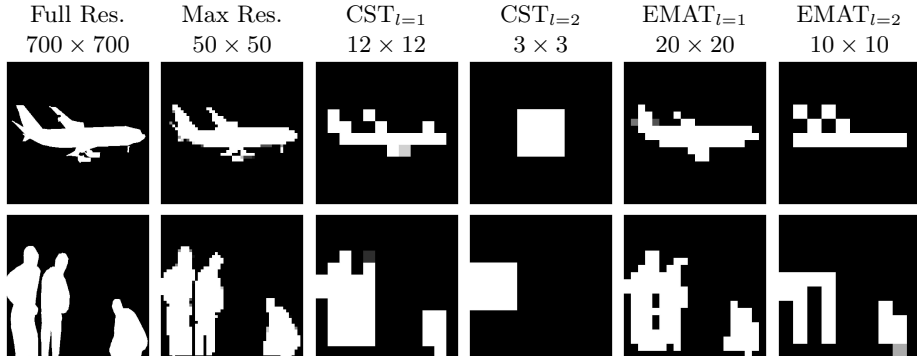


Fig. 5. Segmentation masks used by CST [20] and EMAT. “Full Res.” shows the mask at full resolution, while “Max Res.” is to the highest resolution compatible with the masked-attention layers of both methods. “CST_{*l*= \cdot} ” and “EMAT_{*l*= \cdot} ” indicate the mask resolution used in layer $l \in \{1, 2\}$ of CST and EMAT, respectively.

This ability to preserve fine details explains why EMAT produces more accurate masks especially for small objects, *e.g.*, the dog in the last row of the top part of Fig. 4, the handle of the beer glass in the bottom part of Fig. 4, and the boat and person in Fig. 1.

5.3 Analysis of Small Objects

To further analyze the impact of higher-resolution correlation tokens on small objects, we filter each fold of PASCAL-5ⁱ [36] and COCO-20ⁱ [31] based on object size, creating three splits: objects occupying 0–5%, 5–10%, and 10–15% of the image (see the supplementary material for details on how these splits were defined). Fig. 6 shows the average accuracy and mIoU of CST* and the corresponding improvement achieved by EMAT across the three splits for both datasets. The results indicate that accuracy and mIoU increase with object size, and EMAT provides the largest improvement over CST* for the smallest objects, gradually decreasing as object size increases. The enhanced classification and segmentation accuracy of EMAT is likely due to better localization enabled by its higher-resolution correlation tokens (see Fig. 5).

5.4 Ablation Study

Tab. 2 first reports the results of CST* with its original support dimension per layer t_s^l . For fair comparison, we increased the t_s^l of CST* to use the same as EMAT, but it required about 63 GB of GPU memory, which exceeded the 48 GB capacity of our GPUs, so we instead use the largest t_s^l that fits in our memory. For EMAT we progressively integrated: (1) memory-efficient masked attention (see Sec. 4.1), (2) learnable downscaling of the query matrix (see Sec. 4.2), and (3) parameter-efficiency modifications (see Sec. 4.3). Although Tab. 2 includes

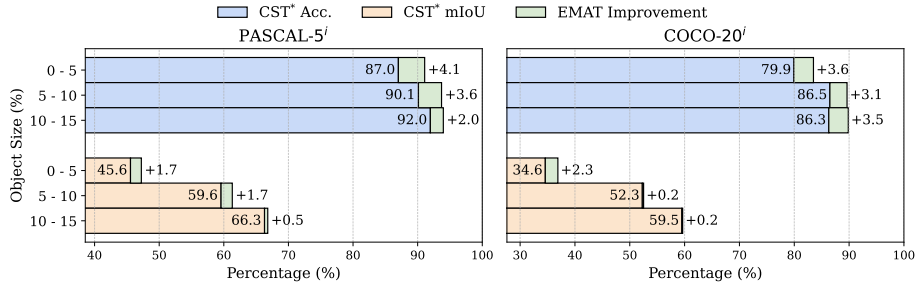


Fig. 6. Analysis of small objects on PASCAL-5ⁱ [36] and COCO-20ⁱ [31]. Each bar represents the average across the four folds of each dataset, filtered by object size, using 1-way 1-shot tasks. To enable a more controlled analysis, we modified the original setting (see Sec. 3) to ensure that the query image always contain the class of the support image. CST* uses the same backbone as EMAT (*i.e.*, DINOv2 [32]).

Table 2. Ablation study of EMAT on PASCAL-5ⁱ [36] under the original evaluation setting. “ t_s^l ” indicates the value of t_s for each layer $l \in \{1, 2\}$. The memory efficiency (ME), learnable downscaling (LD), and parameter efficiency (PE) columns correspond to the modifications described in Secs. 4.1 to 4.3, respectively. “Mem. Usage” reports the average per-GPU memory used during training. “All Dataset” refers to 2-way 1-shot evaluation on the full test set, while “Small Objects” restricts evaluation to objects occupying less than 15% of the image, using 1-way 1-shot tasks in which the query always contains the support class. CST* uses the same backbone as EMAT (*i.e.*, DINOv2 [32]). **(Top)** CST* with its original support dimension per layer t_s^l . **(Middle)** CST* with the largest t_s^l that fits in our 48 GB GPUs. **(Bottom)** successive modifications introduced by EMAT. Highlight indicates our complete model. **Bold** and underlined values indicate the best and second best results.

t_s^l per Layer	Method	ME	LD	PE	Mem. Usage	Train. Params.	All Dataset		Small Objects	
							Acc.	mIoU	Acc.	mIoU
$t_s^1=145$ $t_s^2=10$	CST*	–	–	–	8.68	<u>366.00</u>	80.58	63.28	88.96	58.16
$t_s^1=325$ $t_s^2=37$	CST*	–	–	–	39.22	<u>366.00</u>	<u>82.23</u>	63.31	89.65	<u>58.75</u>
	CST*	–	–	–	≈ 63	<u>366.00</u>	N/A	N/A	N/A	N/A
$t_s^1=401$	EMAT	✓	–	–	36.92	<u>366.00</u>	81.95	62.97	87.99	58.06
$t_s^2=101$	EMAT	✓	✓	–	<u>36.53</u>	404.48	82.17	<u>63.36</u>	<u>90.49</u>	58.73
	EMAT	✓	✓	✓	38.31	86.02	82.70	63.38	91.74	59.17

results on the full PASCAL-5ⁱ test set, the discussion below focuses on the small-object subset to highlight the effect of each modification introduced by EMAT.

Adding our memory-efficient masked attention alone lowers memory usage by 26 GB ($\approx 41\%$) but does not improve accuracy or mIoU compared to either variant of CST*, likely because the model relies on large pooling windows for processing the higher-resolution correlation tokens. Incorporating our learnable

downscaling removes those large windows and yields absolute accuracy gains of +1.53% over the original CST* and of +0.84% over the variant with the larger t_s^l . It also achieves an absolute mIoU gain of +0.57% compared with the original CST*, while matching the mIoU of the variant with larger t_s^l .

Because our learnable downscaling increases the number of trainable parameters, we next apply our parameter-efficiency modifications that remove 318 K parameters ($\approx 79\%$), while still saving about 39% of the memory CST* would need for using the same t_s^l as EMAT. These modifications result in absolute accuracy gains of +2.78% over the original CST* and +2.09% over the variant with the larger t_s^l ; mIoU improves by +1.01% and +0.42%, respectively. EMAT also slightly improves accuracy and mIoU on the full test set, but its largest gains appear on images containing small objects.

6 Limitations

While the results of Tab. 2 validate that our proposed EMAT is memory and parameter efficient for high-resolution correlation tokens, its memory efficiency is constrained to datasets where the segmentation masks contain unlabeled areas. This limitation arises because our memory-efficient masked attention mechanism is equivalent to self-attention [12, 47] when applied to dense semantic-segmentation datasets like Cityscapes [9], where each pixel in the segmentation mask corresponds directly to one of the semantic classes in the image.

7 Conclusion

In this work, we propose EMAT, an enhancement over CST, the state-of-the-art method for few-shot classification and segmentation (FS-CS). EMAT incorporates our novel memory-efficient masked attention mechanism that allows our model to process high-resolution correlation tokens while maintaining memory and parameter efficiency. Our results demonstrate that EMAT consistently outperforms all FS-CS methods across all evaluation settings while requiring at least four times fewer trainable parameters. Moreover, our qualitative results highlight that EMAT is capable of correctly generating empty segmentation masks when necessary and capturing finer details more accurately, which improves accuracy when dealing with small objects. Additionally, we introduce two novel few-shot evaluation settings designed to maximize the use of the available annotations during inference, reflecting practical few-shot scenarios.

Acknowledgments. This work was funded by the Hessian Ministry of Science and Research, Arts and Culture (HMWK) through the project “The Third Wave of Artificial Intelligence – 3AI”. The work was further supported by the Deutsche Forschungsgemeinschaft (German Research Foundation, DFG) under Germany’s Excellence Strategy (EXC 3057/1 “Reasonable Artificial Intelligence”, Project No. 533677015). Stefan Roth acknowledges support by the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement No. 866008).

References

1. Aggarwal, P., Deshpande, A., Narasimhan, K.R.: SemSup-XC: Semantic supervision for zero and few-shot extreme classification. In: ICML. vol. 202, pp. 228–247 (2023) 1
2. Alfassy, A., Karlinsky, L., Aides, A., Shtok, J., Harary, S., Feris, R.S., Giryes, R., Bronstein, A.M.: LaSO: Label-set operations networks for multi-label few-shot learning. In: CVPR. pp. 6548–6557 (2019) 2
3. Allen, K.R., Shelhamer, E., Shin, H., Tenenbaum, J.B.: Infinite mixture prototypes for few-shot learning. In: ICML. vol. 97, pp. 232–241 (2019) 1, 2
4. Antoniou, A., Edwards, H., Storkey, A.J.: How to train your MAML. In: ICLR (2019) 2
5. Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., Joulin, A.: Emerging properties in self-supervised vision transformers. In: ICCV. pp. 9650–9660 (2021) 1
6. Carrión-Ojeda, D., Alam, M., Escalera, S., et al.: NeurIPS’22 Cross-Domain MetaDL Challenge: Results and lessons learned. In: NeurIPS Competition Track. vol. 220, pp. 50–72 (2022) 2
7. Chen, H., Dong, Y., Lu, Z., Yu, Y., Han, J.: Pixel matching network for cross-domain few-shot segmentation. In: WACV. pp. 978–987 (2024) 3, vii
8. Chen, Z., Wu, J., Wang, W., Su, W., Chen, G., Xing, S., Zhong, M., Zhang, Q., Zhu, X., Lu, L., Li, B., Luo, P., Lu, T., Qiao, Y., Dai, J.: InternVL: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In: CVPR. pp. 24185–24198 (2024) 1
9. Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth, S., Schiele, B.: The Cityscapes dataset for semantic urban scene understanding. In: CVPR. pp. 3213–3223 (2016) 13
10. Dhillon, G.S., Chaudhari, P., Ravichandran, A., Soatto, S.: A baseline for few-shot image classification. In: ICLR (2020) 2
11. Dong, N., Xing, E.P.: Few-shot semantic segmentation with prototype learning. In: BMVC (2018) 1, 3
12. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N.: An image is worth 16x16 words: Transformers for image recognition at scale. In: ICLR (2021) 5, 6, 8, 13
13. Fan, X., Wang, X., Gao, J., Wang, J., Luo, Z., Liu, R.: Bi-level learning of task-specific decoders for joint registration and one-shot medical image segmentation. In: CVPR. pp. 11726–11735 (2024) 1, 2
14. Fang, Z., Wang, X., Li, H., Liu, J., Hu, Q., Xiao, J.: FastRecon: Few-shot industrial anomaly detection via fast feature reconstruction. In: ICCV. pp. 17481–17490 (2023) 1
15. Finn, C., Abbeel, P., Levine, S.: Model-agnostic meta-learning for fast adaptation of deep networks. In: ICML. vol. 70, pp. 1126–1135 (2017) 2
16. Gong, X., Xia, X., Zhu, W., Zhang, B., Doermann, D., Zhuo, L.: Deformable Gabor feature networks for biomedical image classification. In: WACV. pp. 4004–4012 (2021) 2
17. Hao, F., He, F., Liu, L., Wu, F., Tao, D., Cheng, J.: Class-aware patch embedding adaptation for few-shot image classification. In: ICCV. pp. 18905–18915 (2023) 2
18. Herzog, J.: Adapt before comparison: A new perspective on cross-domain few-shot segmentation. In: CVPR. pp. 23605–23615 (2024) 1, 2

19. Kang, D., Cho, M.: Integrative few-shot learning for classification and segmentation. In: CVPR. pp. 9979–9990 (2022) 1, 2, 3, 8, 9, i, iv, v
20. Kang, D., Koniusz, P., Cho, M., Murray, N.: Distilling self-supervised vision transformers for weakly-supervised few-shot classification & segmentation. In: CVPR. pp. 19627–19638 (2023) 1, 2, 3, 5, 6, 8, 9, 10, 11, iv, v
21. Kang, D., Kwon, H., Min, J., Cho, M.: Relational embedding for few-shot classification. In: ICCV. pp. 8822–8833 (2021) 2
22. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: ICLR (2015) 8
23. Kirillov, A., Mintun, E., Ravi, N., Mao, H., Rolland, C., Gustafson, L., Xiao, T., Whitehead, S., Berg, A.C., Lo, W.Y., Dollar, P., Girshick, R.: Segment anything. In: ICCV. pp. 4015–4026 (2023) 1
24. Li, W.H., Liu, X., Bilen, H.: Cross-domain few-shot learning with task-specific adapters. In: CVPR. pp. 7161–7170 (2022) 2
25. Liu, H., Li, C., Wu, Q., Lee, Y.J.: Visual instruction tuning. In: NeurIPS (2023) 1
26. Liu, J., Bao, Y., Xie, G.S., Xiong, H., Sonke, J.J., Gavves, E.: Dynamic prototype convolution network for few-shot semantic segmentation. In: CVPR. pp. 11553–11562 (2022) 3
27. Liu, Y., Zhu, M., Li, H., Chen, H., Wang, X., Shen, C.: Matcher: Segment anything with one shot using all-purpose feature matching. In: ICLR (2024) 3, vii
28. Ma, T., Sun, Y., Yang, Z., Yang, Y.: ProD: Prompting-to-disentangle domain knowledge for cross-domain few-shot image classification. In: CVPR. pp. 19754–19763 (2023) 2
29. Min, J., Kang, D., Cho, M.: Hypercorrelation squeeze for few-shot segmentation. In: ICCV. pp. 6941–6952 (2021) 3, 9, iv, v
30. Moon, S., Sohn, S.S., Zhou, H., Yoon, S., Pavlovic, V., Khan, M.H., Kapadia, M.: MSI: Maximize support-set information for few-shot segmentation. In: ICCV. pp. 19266–19276 (2023) 1, 3, vii
31. Nguyen, K., Todorovic, S.: Feature weighting and boosting for few-shot segmentation. In: ICCV. pp. 622–631 (2019) 8, 10, 11, 12, ii, iii, v, vi, vii
32. Oquab, M., Darcet, T., Moutakanni, T., et al.: DINOv2: Learning robust visual features without supervision. *Trans. Mach. Learn. Res.* (2024) 1, 2, 8, 9, 12, iv, v, vi
33. Peng, B., Tian, Z., Wu, X., Wang, C., Liu, S., Su, J., Jia, J.: Hierarchical dense correlation distillation for few-shot segmentation. In: CVPR. pp. 23641–23651 (2023) 3, vii
34. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., Sutskever, I.: Learning transferable visual models from natural language supervision. In: ICML. vol. 139, pp. 8748–8763 (2021) 1
35. Raghu, A., Raghu, M., Bengio, S., Vinyals, O.: Rapid learning or feature reuse? Towards understanding the effectiveness of MAML. In: ICLR (2020) 2
36. Shaban, A., Bansal, S., Liu, Z., Essa, I., Boots, B.: One-shot learning for semantic segmentation. In: BMVC (2017) 8, 11, 12, ii, iv, vi, vii
37. Silva-Rodríguez, J., Hajimiri, S., Ben Ayed, I., Dolz, J.: A closer look at the few-shot adaptation of large vision-language models. In: CVPR. pp. 23681–23690 (2024) 2
38. Simon, C., Koniusz, P., Harandi, M.: Meta-learning for multi-label few-shot classification. In: WACV. pp. 346–355 (2022) 2
39. Snell, J., Swersky, K., Zemel, R.S.: Prototypical networks for few-shot learning. In: NIPS. pp. 4077–4087 (2017) 2

40. Sun, Q., Liu, Y., Chua, T.S., Schiele, B.: Meta-transfer learning for few-shot learning. In: CVPR. pp. 403–412 (2019) 1, 2
41. Team, G.G.: Gemini: A family of highly capable multimodal models. arXiv:2312.11805 [cs.CL] (2023) 1
42. Tian, P., Wu, Z., Qi, L., Wang, L., Shi, Y., Gao, Y.: Differentiable meta-learning model for few-shot semantic segmentation. In: AAAI. pp. 12087–12094 (2020) 3
43. Tian, Y., Wang, Y., Krishnan, D., Tenenbaum, J.B., Isola, P.: Rethinking few-shot image classification: A good embedding is all you need? In: ECCV. vol. 12359, pp. 266–282 (2020) 1, 2
44. Tian, Z., Lai, X., Jiang, L., Liu, S., Shu, M., Zhao, H., Jia, J.: Generalized few-shot semantic segmentation. In: CVPR. pp. 11563–11572 (2022) 1, 5
45. Tian, Z., Zhao, H., Shu, M., Yang, Z., Li, R., Jia, J.: Prior guided feature enrichment network for few-shot segmentation. IEEE T. Pattern Anal. Mach. Intell. **44**(2), 1050–1065 (2022) 3, 9, iv, v
46. Ullah, I., Carrión-Ojeda, D., Escalera, S., Guyon, I., Huisman, M., Mohr, F., van Rijn, J.N., Sun, H., Vanschoren, J., Vu, P.A.: Meta-Album: Multi-domain meta-dataset for few-shot image classification. In: NeurIPS. vol. 35, pp. 3232–3247 (2022) 2
47. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. In: NIPS. pp. 5998–6008 (2017) 6, 13
48. Vinyals, O., Blundell, C., Lillicrap, T., Kavukcuoglu, K., Wierstra, D.: Matching networks for one shot learning. In: NIPS. pp. 3630–3638 (2016) 2, 3, i
49. Wang, J., Zhang, B., Pang, J., Chen, H., Liu, W.: Rethinking prior information generation with CLIP for few-shot segmentation. In: CVPR. pp. 3941–3951 (2024) 3, vii
50. Wang, K., Liew, J.H., Zou, Y., Zhou, D., Feng, J.: PANet: Few-shot image semantic segmentation with prototype alignment. In: ICCV. pp. 9196–9205 (2019) 3, 9, iv, v
51. Wang, Y., Luo, N., Zhang, T.: Focus on query: Adversarial mining transformer for few-shot segmentation. In: NeurIPS. vol. 36, pp. 31524–31542 (2023) 3, vii
52. Wu, A., Han, Y., Zhu, L., Yang, Y.: Universal-prototype enhancing for few-shot object detection. In: ICCV. pp. 9567–9576 (2021) 1
53. Xie, G.S., Xiong, H., Liu, J., Yao, Y., Shao, L.: Few-shot semantic segmentation with cyclic memory network. In: ICCV. pp. 7293–7302 (2021) 3
54. Xu, Q., Zhao, W., Lin, G., Long, C.: Self-calibrated cross attention network for few-shot segmentation. In: ICCV. pp. 655–665 (2023) 3, vii
55. Yang, Y., Chen, Q., Feng, Y., Huang, T.: MIANet: Aggregating unbiased instance and general information for few-shot semantic segmentation. In: CVPR. pp. 7131–7140 (2023) 1, vii
56. Ye, C., Zhu, H., Liao, Y., Zhang, Y., Chen, T., Fan, J.: What makes for effective few-shot point cloud classification? In: WACV. pp. 1829–1838 (2022) 1
57. Zhang, A., Gao, G., Jiao, J., Liu, C., Wei, Y.: Bridge the points: Graph-based few-shot segment anything semantically. In: NeurIPS (2024) 3, vii
58. Zhang, B., Xiao, J., Qin, T.: Self-guided and cross-guided learning for few-shot segmentation. In: CVPR. pp. 8312–8321 (2021) 3
59. Zhang, M., Shi, M., Li, L.: MFNet: Multiclass few-shot segmentation network with pixel-wise metric learning. IEEE T. Circuits Syst. Video Tech. **32**(12), 8586–8598 (2022) 3
60. Zhou, F., Wang, P., Zhang, L., Wei, W., Zhang, Y.: Revisiting prototypical network for cross domain few-shot learning. In: CVPR. pp. 20061–20070 (2023) 2

61. Zhou, Z., Xu, H.M., Shu, Y., Liu, L.: Unlocking the potential of pre-trained vision transformers for few-shot semantic segmentation through relationship descriptors. In: CVPR. pp. 3817–3827 (2024) 3
62. Zhu, L., Chen, T., Ji, D., Ye, J., Liu, J.: LLaFS: When large language models meet few-shot segmentation. In: CVPR. pp. 3065–3075 (2024) 1, 3, vii
63. Zhu, L., Chen, T., Yin, J., See, S., Liu, J.: Addressing background context bias in few-shot segmentation through iterative modulation. In: CVPR. pp. 3370–3379 (2024) 1, vii
64. Zhu, X., Zhang, R., He, B., Zhou, A., Wang, D., Zhao, B., Gao, P.: Not all features matter: Enhancing few-shot clip with adaptive prior refinement. In: ICCV. pp. 2605–2615 (2023) 2
65. Zhu, Y., Wang, S., Xin, T., Zhang, H.: Few-shot medical image segmentation via a region-enhanced prototypical transformer. In: MICCAI. pp. 271–280 (2023) 1, 2

A Few-shot Evaluation Settings

As explained in Sec. 3.1, we introduce two novel few-shot evaluation settings, partially and fully augmented, to better utilize available annotations and represent realistic scenarios. Fig. 7 illustrates how the support set of a 2-way 2-shot (base task configuration) task is modified under each setting.

In the *original setting* [19], additional annotations are discarded, leaving the base task configuration unchanged. The drawback of this setting is that it ignores available annotations. For example, in Fig. 7, the second image for the 1st class (dog) also contains annotations for the 2nd class (horse), but these are removed.

To address the loss of available annotations, our *partially augmented setting* retains annotations for the support classes (*e.g.*, dog and horse in Fig. 7) while removing those for non-support classes (*e.g.*, person in Fig. 7). Since few-shot classification and segmentation (FS-CS) methods treat N -way K -shot tasks as N separate 1-way K -shot tasks, each image (shot) must contain only one class (way). Therefore, if an image contains multiple support classes, we duplicate it so each copy includes only one. This leads to a task configuration that no longer follows the N -way K -shot definition [48], as the number of shots per class can vary. Consequently, the classes with fewer shots are randomly augmented using the available support data. For instance, in Fig. 7, the partially augmented setting converts the 2-way 2-shot into a 2-way 3-shot task. It is important to note that in this setting, the number of ways remains fixed but the number of shots can increase up to $N \times K$.

On the other hand, our *fully augmented setting* incorporates all available annotations, including those for non-support classes (*e.g.*, person in Fig. 7). The task is then processed similarly to the partially augmented setting. However, in this case, both the number of ways and the number of shots can increase, with shots still bounded by $N \times K$. Regardless of the setting, models are evaluated only on the final support classes. For example, in Fig. 7, models in the original and partially augmented settings are evaluated on dogs and horses, while in the fully augmented setting, they are also evaluated on people.

Practical Considerations. Because annotating data is time- and resource-consuming, few-shot evaluation settings should aim to maximize the use of existing annotations. Our partially augmented setting achieves this without increasing task difficulty and may even enhance FS-CS performance by providing more examples per class. On the other hand, our fully augmented setting incorporates all available annotations, increasing both the number of ways and the number of shots, making the task more challenging by requiring the model to classify and segment additional classes. Nevertheless, neither setting requires extra annotations beyond those already available; instead, they use existing labels discarded by the original setting to better represent realistic scenarios.

Although the loss of annotations in the original setting may not significantly affect training, our partially augmented setting can safely be used in that phase. In contrast, our fully augmented setting should be used only for evaluation to

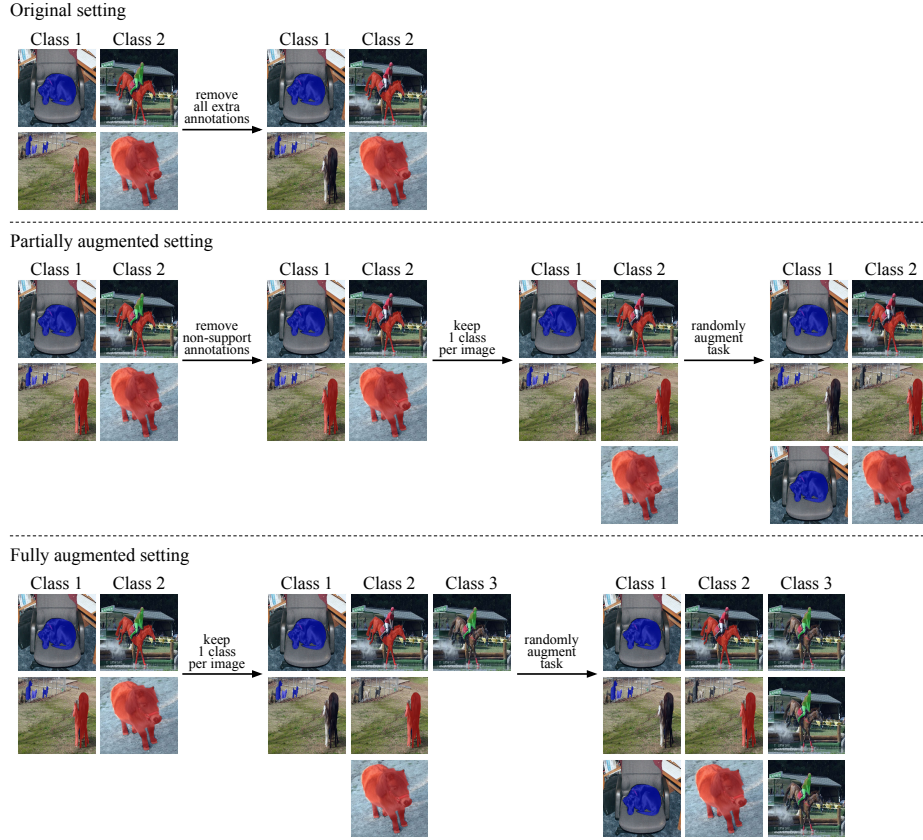


Fig. 7. Task augmentation pipeline for each few-shot evaluation setting. The 2-way 2-shot (base configuration) task is extracted from the PASCAL-5ⁱ dataset [36]. The *original setting* leaves this configuration unchanged, discarding available annotations. The *partially augmented setting* incorporates unused support-class annotations, yielding a 2-way 3-shot task. The *fully augmented setting* includes all annotations, increasing complexity to 3-way 3-shot.

prevent data leakage, as it mixes training and test classes within the support set. For this reason, we recommend using both settings exclusively for evaluation.

Object Size Filtering. Very small objects can add noise during task augmentation in the partially and fully augmented settings, so we filter out any object that occupy less than a threshold θ of the image area. We set $\theta = 7\%$ for PASCAL-5ⁱ [36] and $\theta = 3\%$ COCO-20ⁱ [31], values that correspond to the 25th percentile of object sizes in the training set of each dataset. These different thresholds reflect dataset-specific object size distribution.

Fig. 8 shows the impact of applying these thresholds on the classification accuracy and segmentation mIoU of our efficient masked attention transformer

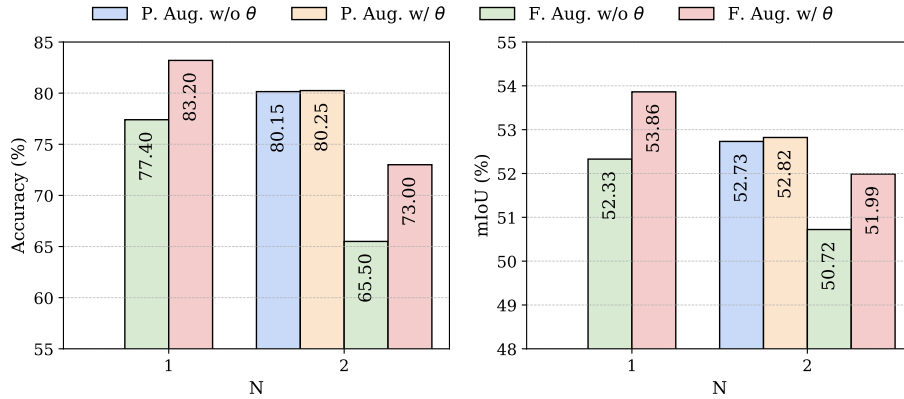


Fig. 8. Impact of object size filtering on task augmentation. “P. Aug.” and “F. Aug.” denote the partially and fully augmented settings, respectively, while θ indicates whether object size filtering was applied. The bars show the accuracy or mIoU of EMAT on the COCO-20ⁱ dataset [31] for {1, 2}-way 1-shot tasks. Results for the 1-way 1-shot configuration in the partially augmented setting are omitted because, when $N = 1$, this setting is equivalent to the original one. In both augmented settings, removing small objects increases accuracy and mIoU.

(EMAT) on COCO-20ⁱ for {1, 2}-way 1-shot tasks. Filtering out small objects consistently improves accuracy and mIoU in both augmented settings, with a larger gain in the fully augmented setting because it augments more tasks. Therefore, all results were obtained using the proposed thresholding.

B FS-CS Per-Fold Results

Sec. 5.1 compares our EMAT with current state-of-the-art (SOTA) FS-CS methods, reporting results averaged over the four folds of the PASCAL-5ⁱ and COCO-20ⁱ datasets using 2-way 1-shot tasks. To complement those results, Tabs. 3 and 4 present per-fold results, consistently showing that EMAT outperforms all FS-CS methods across most folds, datasets, and evaluation settings. These tables also report the number of augmented tasks in our proposed evaluation settings. On PASCAL-5ⁱ, the partially and fully augmented settings augment 242 and 694 tasks, respectively, meaning that about 6 % of tasks contain extra support-class annotations (partially augmented), and 11 % contain extra annotations for non-support classes (fully augmented). On COCO-20ⁱ, the corresponding augmentation rates are roughly 3 % and 35 %. These results confirm that our evaluation settings use annotations discarded by the original setting, thereby creating more realistic evaluation scenarios.

To evaluate the scalability across different N -way K -shot configurations, Fig. 9 shows the accuracy and mIoU of all FS-CS methods on PASCAL-5ⁱ and COCO-20ⁱ in the fully augmented setting using {1, ..., 5}-way 5-shot tasks. EMAT consistently achieves higher classification accuracy and mIoU compared

Table 3. Per-fold comparison of FS-CS methods on PASCAL-5ⁱ [36]. Results are reported for all evaluation settings (original, partially augmented, and fully augmented) using 1000 tasks per fold with a base configuration of 2-way 1-shot. The number in parentheses in the partially and fully augmented settings indicates the total number of augmented tasks. CST* and EMAT were trained and evaluated, while other methods were only evaluated using the checkpoints from [19]. CST* uses the same backbone as EMAT (*i.e.*, DINOv2 [32]). All values are percentages (higher is better). Highlight indicates our proposed method. **Bold** and underlined values represent the best and second best results.

Method	Fold-0		Fold-1		Fold-2		Fold-3		Average	
	Acc.	mIoU	Acc.	mIoU	Acc.	mIoU	Acc.	mIoU	Acc.	mIoU
Original										
PANet [50]	62.30	33.31	53.30	45.94	49.50	31.20	61.00	38.34	56.53	37.20
PFENet [45]	23.10	31.48	53.80	46.60	40.60	31.15	39.90	33.03	39.35	35.57
HSNet [29]	68.20	43.97	73.00	55.12	56.90	35.19	71.00	45.14	67.27	44.85
ASNet [19]	68.20	48.44	76.20	58.19	58.80	36.45	70.00	48.41	68.30	47.87
CST [20]	70.10	53.90	75.20	59.98	61.70	46.30	74.50	54.93	70.37	53.78
CST*	<u>89.40</u>	<u>64.61</u>	<u>80.90</u>	68.16	<u>71.40</u>	55.50	<u>80.60</u>	<u>64.84</u>	<u>80.58</u>	<u>63.28</u>
EMAT	90.70	66.68	83.50	<u>67.61</u>	72.00	<u>54.15</u>	84.60	65.07	82.70	63.38
Partially Augmented (↑ 242 tasks)										
PANet [50]	62.30	33.31	53.00	46.00	51.70	32.35	60.70	38.32	56.93	37.49
PFENet [45]	23.10	31.48	54.10	46.68	40.70	31.28	40.00	33.02	39.48	35.61
HSNet [29]	68.20	43.97	72.90	54.85	59.00	34.90	70.90	45.18	67.75	44.72
ASNet [19]	68.20	48.44	76.00	57.88	60.50	36.28	69.80	48.51	68.62	47.78
CST [20]	70.10	53.90	75.30	60.04	62.40	46.32	74.60	54.97	70.60	53.81
CST*	<u>89.40</u>	<u>64.61</u>	<u>80.90</u>	68.10	<u>71.50</u>	55.44	<u>80.60</u>	<u>64.77</u>	<u>80.60</u>	<u>63.23</u>
EMAT	90.70	66.68	83.30	<u>67.48</u>	73.00	<u>54.03</u>	84.70	65.10	82.92	63.32
Fully Augmented (↑ 694 tasks)										
PANet [50]	62.30	33.31	51.70	45.95	49.60	31.65	59.40	38.11	55.75	37.25
PFENet [45]	23.10	31.48	52.60	46.14	33.70	30.33	38.10	32.35	36.88	35.08
HSNet [29]	68.20	43.97	72.20	54.55	53.80	33.87	69.50	45.19	65.92	44.40
ASNet [19]	68.20	48.44	75.00	57.66	54.30	35.90	68.10	48.31	66.40	47.58
CST [20]	70.10	53.90	74.50	59.50	56.70	46.62	72.50	55.01	68.45	53.76
CST*	<u>89.40</u>	<u>64.61</u>	<u>80.20</u>	67.43	<u>66.40</u>	55.80	<u>78.30</u>	<u>64.47</u>	<u>78.57</u>	<u>63.08</u>
EMAT	90.70	66.68	82.60	<u>66.93</u>	68.30	<u>54.37</u>	83.30	65.00	81.23	63.24

to other FS-CS methods. Notably, while increasing the number of classes (N) typically raises task difficulty, EMAT maintains stable segmentation performance, further validating the fully augmented setting as a challenging yet realistic and effective benchmark for FS-CS evaluation.

C Analysis of Object Size Distribution in PASCAL-5ⁱ

Sec. 5.3 analyzes the impact of higher-resolution correlation tokens on small objects by filtering each fold of PASCAL-5ⁱ and COCO-20ⁱ based on object

Table 4. Per-fold comparison of FS-CS methods on COCO-20ⁱ [31]. Results are reported for all evaluation settings (original, partially augmented, and fully augmented) using 1000 tasks per fold with a base configuration of 2-way 1-shot. The number in parentheses in the partially and fully augmented settings indicates the total number of augmented tasks. CST* and EMAT were trained and evaluated, while other methods were only evaluated using the checkpoints from [19]. CST* uses the same backbone as EMAT (*i.e.*, DINOv2 [32]). All values are percentages (higher is better). **Highlight** indicates our proposed method. **Bold** and underlined values represent the best and second best results.

Method	Fold-0		Fold-1		Fold-2		Fold-3		Average	
	Acc.	mIoU	Acc.	mIoU	Acc.	mIoU	Acc.	mIoU	Acc.	mIoU
Original										
PANet [50]	46.60	24.91	52.70	24.98	55.90	23.31	50.00	21.36	51.30	23.64
PFENet [45]	35.60	23.99	34.30	24.57	43.10	20.99	32.80	23.93	36.45	23.37
HSNet [29]	57.70	29.77	62.30	30.94	67.10	31.31	62.60	30.31	62.43	30.58
ASNet [19]	<u>59.50</u>	29.75	61.50	32.99	68.80	33.41	62.40	30.35	63.05	31.62
CST [20]	61.00	34.74	66.40	37.14	68.20	36.76	60.50	36.29	64.02	36.23
CST*	74.00	<u>49.38</u>	<u>79.90</u>	<u>53.88</u>	<u>81.30</u>	<u>51.46</u>	<u>79.60</u>	<u>51.15</u>	<u>78.70</u>	<u>51.47</u>
EMAT	74.00	50.54	83.10	55.44	83.10	53.05	80.10	52.19	80.07	52.81
Partially Augmented (↑ 106 tasks)										
PANet [50]	46.40	25.27	52.80	25.10	55.90	23.27	50.20	21.48	51.32	23.78
PFENet [45]	35.80	24.08	34.30	24.52	43.00	21.00	32.90	23.96	36.50	23.39
HSNet [29]	57.50	29.88	62.60	31.02	67.00	31.30	62.50	30.44	62.40	30.66
ASNet [19]	59.30	29.87	61.70	32.95	68.80	33.39	62.30	30.36	63.03	31.64
CST [20]	61.30	34.59	66.20	37.08	68.40	36.76	60.50	36.36	64.10	36.20
CST*	74.60	<u>49.27</u>	<u>79.80</u>	<u>53.98</u>	<u>81.50</u>	<u>51.68</u>	<u>79.60</u>	<u>51.19</u>	<u>78.87</u>	<u>51.53</u>
EMAT	<u>74.30</u>	50.59	83.20	55.42	83.30	53.06	80.20	52.22	80.25	52.82
Fully Augmented (↑ 1515 tasks)										
PANet [50]	30.90	24.43	49.30	24.44	53.50	22.98	46.60	20.81	45.07	23.17
PFENet [45]	19.60	20.33	29.10	23.79	40.70	20.40	27.90	21.92	29.33	21.61
HSNet [29]	40.20	27.73	57.50	29.82	64.30	30.96	<u>58.60</u>	29.26	55.15	29.44
ASNet [19]	41.40	27.23	55.70	32.10	66.60	33.14	58.20	29.43	55.47	30.47
CST [20]	44.70	33.74	59.70	36.32	65.60	36.87	55.20	35.46	56.30	35.60
CST*	<u>54.40</u>	<u>48.25</u>	<u>74.80</u>	<u>52.86</u>	<u>79.70</u>	<u>51.69</u>	75.80	<u>50.25</u>	<u>71.18</u>	<u>50.76</u>
EMAT	56.30	49.53	78.50	54.46	81.40	53.06	75.80	50.90	73.00	51.99

size. This filtering results in three splits: objects occupying 0–5 %, 5–10 %, and 10–15 % of the image. These intervals were chosen after examining the object size distribution in the PASCAL-5ⁱ test set, which has fewer test examples than COCO-20ⁱ.

Fig. 10 presents frequency histograms of object sizes (up to 50 % of the image) for all four PASCAL-5ⁱ test folds. Based on the fold with more examples (Fold 2), we selected the size intervals that ensure each split contains at least 100 examples, *i.e.*, 0–5 %, 5–10 %, and 10–15 %. For consistency and comparability,

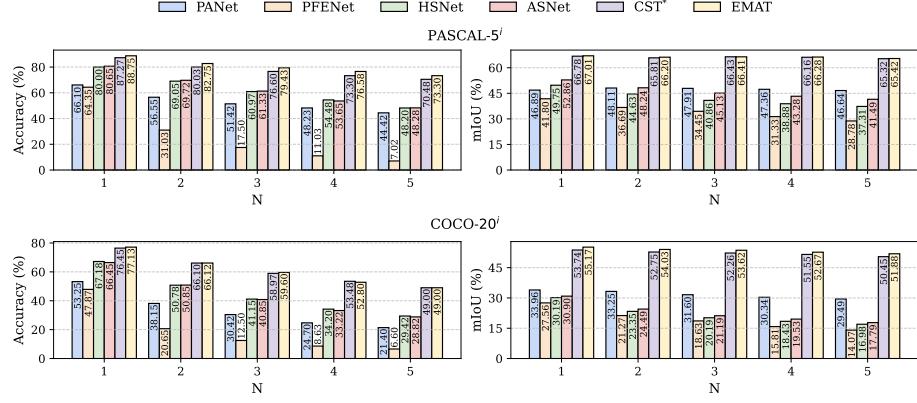


Fig. 9. Comparison of FS-CS methods on PASCAL-5ⁱ [36] and COCO-20ⁱ [31] under the fully augmented setting. CST* uses the same backbone as EMAT (*i.e.*, DINOv2 [32]). Results are averaged over 4000 tasks (1000 per fold) using a base task configuration of {1, ..., 5}-way 5-shot. The number of augmented tasks for each value of N in PASCAL-5ⁱ is 1243, 1758, 2094, 2319, and 2522, respectively. For COCO-20ⁱ, the corresponding values are 2136, 2982, 3432, 3674, and 3806.

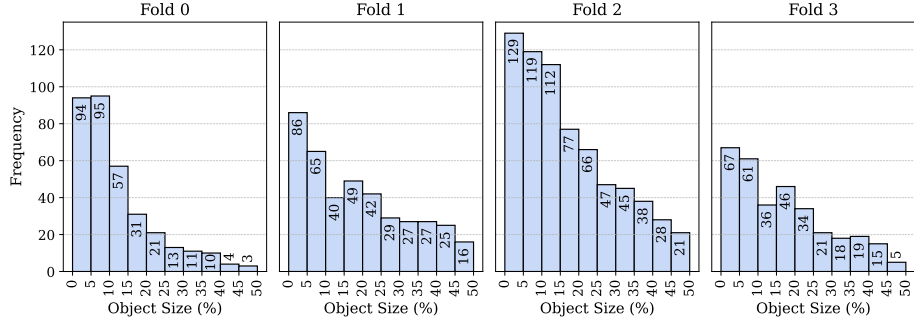


Fig. 10. Per-fold object size distribution in the PASCAL-5ⁱ [36] test set. Each histogram shows the frequency of objects based on their relative size (as a percentage of the image area), up to 50%.

we apply the same splits to COCO-20ⁱ. Furthermore, Fig. 10 shows that the majority of objects in PASCAL-5ⁱ are very small, highlighting the importance of accurately classifying and segmenting small objects, a challenge that our EMAT effectively addresses.

D Comparison to SOTA FS-S

To compare our EMAT with SOTA few-shot segmentation (FS-S) methods, we removed its classification head and denote this adapted version as EMAT_s. This adaptation ensures a fair comparison to existing FS-S methods, as simultane-

Table 5. Comparison of FS-S methods on PASCAL-5ⁱ [36] and COCO-20ⁱ [31] using 1-way {1, 5}-shot tasks. For fairness, EMAT_s exclude the classification head. All values correspond to mIoU in percentage (higher is better). The mIoU values for EMAT_s were obtained through our own training and evaluation, while values for other methods are taken directly from their respective papers. The top part of the table shows methods based on ResNet backbones, and the bottom part shows those based on foundation models. **Highlight** indicates our proposed method. **Bold** and underlined values represent the best and second best results.

Method	Venue	Backbone	PASCAL-5 ⁱ		COCO-20 ⁱ	
			1-s	5-s	1-s	5-s
MIANet [55]	CVPR’23	ResNet50	68.7	71.6	47.7	51.7
HDMNet [33]	CVPR’23	ResNet50	69.4	71.8	50.0	56.0
VAT + MSI [30]	ICCV’23	ResNet101	70.1	72.2	49.8	54.0
SCCAN [54]	ICCV’23	ResNet101	68.3	71.5	48.2	57.0
AMFomer [51]	NeurIPS’23	ResNet101	70.7	73.6	51.0	57.3
PMNet [7]	WACV’24	ResNet101	68.1	73.9	43.7	53.1
ABCB [63]	CVPR’24	ResNet101	72.0	74.9	51.5	58.8
EMAT _s	–	DINOv2-S	72.5	75.9	59.8	65.0
LLaFS [62]	CVPR’24	ResNet50 + CodeLlama-7B	<u>73.5</u>	75.6	53.9	60.0
PI-CLIP [49]	CVPR’24	ResNet50 + CLIP-B	76.8	<u>77.2</u>	56.8	59.1
Matcher [27]	ICLR’24	DINOv2-L + SAH	68.1	74.0	52.7	60.7
GF-SAM [57]	NeurIPS’24	DINOv2-L + SAH	72.1	82.6	<u>58.7</u>	66.8
EMAT _s	–	DINOv2-S	72.5	75.9	59.8	<u>65.0</u>

ously handling classification and segmentation is more challenging than focusing only on segmentation. Furthermore, we adopt the 1-way K -shot formulation used by the FS-S methods, where the query image always contains the support class.

Tab. 5 presents a comparison between EMAT_s and SOTA FS-S methods, divided into two groups: methods based on ResNet backbones and those based on foundation models. Although EMAT_s is based on a foundation model (DINOv2), it uses the small variant (DINOv2-S), whose model size is comparable to ResNet50. Therefore, comparing EMAT_s with FS-S methods based on ResNet backbones is still meaningful.

The top part of Tab. 5 shows that EMAT_s consistently outperforms all ResNet-based FS-S methods across both PASCAL-5ⁱ and COCO-20ⁱ datasets and across all task configurations, even surpassing methods based on a significantly larger backbone, ResNet101. Notably, the improvement in mIoU is more pronounced on COCO-20ⁱ, with absolute gains of +8.3% and +6.2% for 1-shot and 5-shot settings, respectively.

On the other hand, comparing EMAT_s to FS-S methods based on foundation models (bottom part of Tab. 5) is less straightforward for two main reasons: (1) all baseline methods combine two models, and (2) the foundation models they use are significantly larger than the one used by EMAT_s. For example, both Matcher [27] and GF-SAM [57] use the large version of DINOv2 (*i.e.*, DINOv2-

L), while EMAT_s uses the small variant (*i.e.*, DINOv2-S). Nevertheless, our method still achieves an absolute improvement of +1.1 % mIoU on COCO-20ⁱ in the 1-way 1-shot setting and overall obtains competitive mIoU compared to other SOTA FS-S methods based on large foundation models.