

# Distilling Self-Supervised Vision Transformers for Weakly-Supervised Few-Shot Classification & Segmentation – Supplementary Material –

Dahyun Kang<sup>1,2\*</sup> Piotr Koniusz<sup>3,4</sup> Minsu Cho<sup>2</sup> Naila Murray<sup>1</sup>  
<sup>1</sup>Meta AI <sup>2</sup>POSTECH <sup>3</sup>Data61♥CSIRO <sup>4</sup>Australian National University

We provide additional details and analyses of the proposed method in this supplementary material.

## 1. Implementation details

Our framework is implemented using the PyTorch Lightning library [5] and we use the public implementation and pre-trained model checkpoints of DINO [2].<sup>1</sup> All the experiments with DINO ViTs in the main paper use DINO ViT-small with a patch size of  $8 \times 8$ . Input images to the model are resized to  $400 \times 400$  without any data augmentation schemes, and are fed to a ViT-small of  $8 \times 8$  patch size, which returns  $50^2$  patch tokens. We resize the support image token map dimension, originally  $H_s \times W_s$ , to  $12 \times 12$  via bilinear interpolation to reduce memory footprint. The architecture details of CST are illustrated in Fig. 1 and enumerated in Table 1. Pascal-5<sup>i</sup> and COCO-20<sup>i</sup> are derived from Pascal Visual Object Classes 2012 [4] and Microsoft Common Object in Context 2014 [10], respectively. All experiments use four NVIDIA Tesla V100 GPUs.

## 2. Further analyses

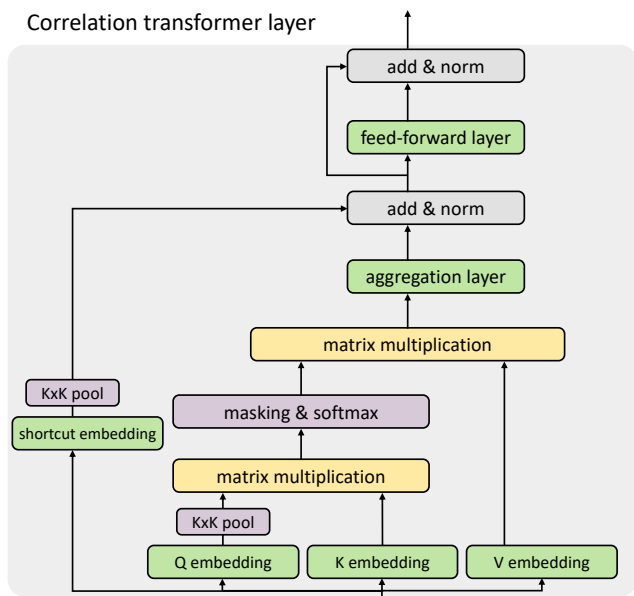
In this section we provide supplementary analyses and results omitted in the main paper due to the page limit.

**Computational complexity.** Table 2 compares the numbers of parameters, MACs, GPU memory consumption of different methods. The computational complexity of each model is evaluated on forwarding a 1-way 1-shot episode of images with  $400 \times 400$  size.

**Effect of self- vs. class-supervised ViT backbone.** Table 3 shows the superior performance of the self-supervised ViT backbone when compared with the class-supervised one, on the task of FS-CS. The class-supervised ViT has the same architecture as the self-supervised one, but was trained on the 1000-class classification task using class supervision on ImageNet 1K [16]. The gap is especially significant

\*Work done during an internship at FAIR.

<sup>1</sup><https://github.com/facebookresearch/dino>.



**Figure 1.** Illustration of correlation transformer layer. Each query-key-value, shortcut, aggregation, and feed-forward layer is implemented with an MLP layer. We use group normalization [21] with 4 groups and ReLU [13] activation in implementation but omit them in this illustration for simplicity.

name	$C_{in} \rightarrow C_{out}$	component	count	# params
correlation transformer	$72 \rightarrow 128$	correlation transformer layers	2	77.5 K
clf. head	$128 \rightarrow 2$	$1 \times 1$ convolutions	2	29.1 K
seg. head	$128 \rightarrow 2$	$3 \times 3$ convolutions	2	259.5 K
total				366.0 K

**Table 1.** Model components of Classification-Segmentation Transformer (CST). Input and output channel dimensions are denoted with  $C_{in}$  and  $C_{out}$ . The backbone network is omitted.

when training our CST model using only image-level labels (33.2% vs. 6.9%). The class-supervised ViTs localize foreground regions less accurately than the self-supervised one (cf. Fig. 3), leading to less accurate pseudo-labels and, ultimately, lower segmentation performance.

model	frozen backbone			training module		GPU
	name	MACs	# params	MACs	# params	memory
HSNet [12]	ResNet50	13.4 G	23.6 M	17.7 G	2.6 M	2.1 G
ASNet [9]	ResNet50	13.4 G	23.6 M	7.3 G	1.3 M	2.3 G
CST (ours)	ViT-S/8	53.4 G	21.7 M	3.7 G	0.4 M	2.4 G

**Table 2.** Comparing computational complexity of different models for forwarding a 1-way 1-shot episode.

ViT backbone	superv.	clf. ER	seg. mIoU
class-sup.	image	79.6	6.9
self-sup. (CST)	image	<b>79.9</b>	<b>33.2</b>
class-sup.	pixel	85.8	54.0
self-sup. (CST)	pixel	<b>85.7</b>	<b>55.5</b>

**Table 3.** Comparing class- and self-supervised ViT backbones for FS-CS on Pascal-5<sup>i</sup> [17].

method	backbone	task heads	clf. ER	seg. mIoU	learn. params.
ASNet [9]	ResNet50	coupled	84.9	52.3	1.4M
CST-(a)	ResNet50	coupled	83.9	51.0	0.4M
CST-(b)	ResNet50	decoupled	84.1	50.9	0.4M
CST-(c)	DINO ViT	coupled	84.3	54.2	0.4M
CST	DINO ViT	decoupled	<b>85.7</b>	<b>55.5</b>	0.4M

**Table 4.** Comparing self-supervised DINO ViT and class-supervised ResNet50 for FS-CS with image-level supervision on Pascal-5<sup>i</sup>.

model	image-level		pixel-level	
	clf. ER (%)	seg. mIoU (%)	clf. ER (%)	seg. mIoU (%)
HSNet* [12]	76.4	31.7	82.2	49.5
ASNet* [9]	78.6	32.8	84.7	53.7
CST (ours)	<b>79.9</b>	<b>33.2</b>	<b>85.7</b>	<b>55.5</b>

**Table 5.** Comparing model performance on FS-CS using the DINO ViT backbone on Pascal-5<sup>i</sup>. The methods denoted with \* are reproduced with the same DINO ViT backbone that CST uses. In image-level supervised learning, all methods are trained with the generated pseudo-GT mask labels.

**Comparing DINO ViT-small vs. ResNet50.** Table 4 compares results using ResNet50 [8] or ViT-small as backbones. Note that CST has independent classification and segmentation task heads that take input from class and image tokens respectively. Therefore, it is not trivial to use ResNet50 as a backbone for CST in order to compare to ResNet-based methods [9, 12, 18]. Albeit not an apple-to-apple comparison, we adapt ResNet50 for CST by artificially creating a “class token” for representing global image semantics, by global average-pooling ResNet feature maps. The ResNet-based CSTs (CST-(a) and (b)) achieve similar performance to ASNet [9] with fewer learnable parameters. Comparing CST-(a) and (b), the gain from using decoupled task heads is unclear, unlike with the ViT-based CST. We hypothesize

backbone	pretrained with	clf. ER	seg. mIoU
ViT [3]	MAE [7]	68.5	13.4
ViT [3]	DINO [2]	82.2	33.7

**Table 6.** Comparing self-supervised DINO and MAE backbones for FS-CS with image-level supervision on Pascal-5<sup>i</sup>.

that, because the class and segmentation representations are generated from the same ResNet features, there is less benefit of task-head specialization when compared to using the class and image tokens in ViTs.

Similarly, Table 5 compares results using DINO ViT-small as the backbones. Other methods [9, 12] are reproduced with the DINO ViT backbone such that they take the same correlation tokens with CST. As those methods [9, 12] benefit from fusing the pyramidal ResNet features, they are not perfectly compatible with DINO ViTs.

**Comparison with other self-supervised backbones.** Table 6 compares the performance of our method when using different self-supervised backbones on FS-CS: DINO ResNet50 and Masked Auto-encoder (MAE) ViT [7]. MAE learns to reconstruct some masked input image tokens in an auto-encoder [15, 19] framework, requiring no supervision. We use an MAE-trained model with an architecture that is identical to the DINO-trained ViTs, and that is also publicly available<sup>2</sup>. We observe that MAE ViTs show weaker localization properties when compared to DINO ViTs as qualitatively compared in Fig. 3, resulting in low segmentation performance.

**Effect of multi-head token correlations.** We examine the effectiveness of head-wise correlation tokens in Eq. (2) in Table 7. We split the class and image tokens into  $M$  equal-sized tokens, each with dimensionality  $C/M$ , compute cross-correlation, and then concatenate the  $M$ -head correlations. We observe that the multi-head token correlation is not only effective in performance but also boosting faster convergence.

**Loss-balancing hyperparameter  $\lambda$ .** We choose the loss-balancing hyperparameter  $\lambda$  based on the course-grained hyperparameter search shown in Table 8. The experimental results show that the performance does not differ significantly, implying that it is robust to different values of  $\lambda$ .

**Four-fold performance.** We omit the four-fold performance for a few experiments in the main paper due to space limitations, and thus specify the four-fold performance of Tables 2, 5, and 8, and Fig. 5 of the main paper in Tables 9, 12, 13, 14, respectively, for reference.

**Experiments with  $K > 1$  shots.** Tables 10 and 11 present FS-CS performance when 5-shot support examples are used

<sup>2</sup><https://github.com/facebookresearch/mae>. As the available models are pre-trained with 16×16 patches, we compare it with DINO ViTs pretrained with the same patch size.

for each class during testing. Using five shots per class brings 3.0% segmentation improvement compared to one-shot models, when zero ground-truth pixel-level labels are used. It is noteworthy that using more shots leads to greater performance gains when learning with pseudo-GT masks, when compared to learning without it; CST with pseudo-GT masks ( $\checkmark$ ) improves from 33.2%  $\rightarrow$  36.2% with 1  $\rightarrow$  5 shots, CST without pseudo-GT masks ( $\times$ ) improves from 16.0%  $\rightarrow$  16.9% (*cf.* Table 1 of the main paper).

**Visualization of classification and segmentation token maps.** Figure 2 visualizes more samples of the channel-averaged classification and segmentation token maps,  $C_{\text{clf}}$  and  $C_{\text{seg}}$ , some of which were also visualized in Fig. 7 of the main paper. The segmentation token maps delineate object boundaries and the foreground, while suppressing background regions.

**Visualization of pseudo-GT masks generated from different pre-trained feature extractors.** Figure 3 visualizes the pseudo-GT masks generated from different models. Pseudo-GT masks from three ViTs with different training procedures are produced as formulated in equations 9 and 10 of the main paper. To generate pseudo-labels from a self-supervised ResNet, the globally-average pooled feature map from the support image is correlated against an image feature map. In addition, we qualitatively observe that, with the ResNet backbone, the inverse correlation better captures foreground objects, and visualize this instead. The qualitative results align with the quantitative results (Table 6 in the main paper), in that the DINO ViT backbone produces the highest-quality pseudo-GT masks.

**Visualization of predicted segmentation masks.** Figure 4 visualizes segmentation masks predicted by the image-level supervised CST. The model correctly segments no foreground when support image classes are irrelevant to those of the query (1st row). Note that the model can also segment small objects, and multiple objects. The last two rows show failure cases; the model segments the apple container which is not present in the ground-truth annotation (5th row) and incorrectly segments a car given the support image containing buses, which may result from vehicular class confusion (6th row). Figure 5 visualizes segmentation masks predicted by the pixel-level supervised CST. Leveraging pixel-level ground-truth annotations, the pixel-level supervised model captures small objects at image corners (2nd and 3rd row) and multiple objects from multiple classes (5th row) precisely.

correlation	superv.	clf. ER	seg. mIoU
single-head	image	<b>80.3</b>	33.1
multi-head (CST)	image	79.9	<b>33.2</b>
single-head	pixel	85.5	54.3
multi-head (CST)	pixel	<b>85.7</b>	<b>55.5</b>

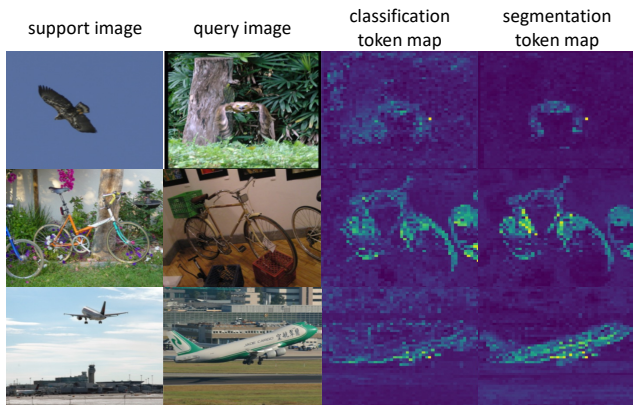
**Table 7.** Comparing single- and multi-head token correlations of Eq. (2) for FS-CS on Pascal-5<sup>i</sup>.

hyperparameter $\lambda$	0.05	0.1	0.5
clf. ER (%)	79.6	79.9	80.5
seg. mIoU (%)	33.1	33.2	32.9

**Table 8.** Hyperparameter search on the loss balancing hyperparameter  $\lambda$  in Eq. (7) in the main paper.

method	N-way 1-shot									
	classification 0/1 exact ratio (%)					segmentation mIoU (%)				
	1	2	3	4	5	1	2	3	4	5
PANet [20]	69.0	50.9	39.3	29.1	22.2	36.2	37.2	37.1	36.6	35.3
PFENet [18]	74.6	41.0	24.9	14.5	7.9	43.0	35.3	30.8	27.6	24.9
HSNet [12]	82.7	67.3	52.5	45.2	36.8	49.7	43.5	39.8	38.1	36.2
ASNet [9]	84.9	68.3	55.8	46.8	<b>37.3</b>	52.3	47.8	45.4	44.5	42.4
CST	<b>85.7</b>	<b>70.4</b>	<b>57.3</b>	<b>47.3</b>	36.9	<b>55.5</b>	<b>53.7</b>	<b>52.6</b>	<b>52.0</b>	<b>50.3</b>

**Table 9.** Numerical performance of Fig. 5 in the main paper: performance comparison on FS-CS with pixel-level supervision on N-way 1-shot Pascal-5<sup>i</sup>.



**Figure 2.** Examples of channel-averaged classification and segmentation task token maps,  $C_{\text{clf}}$  and  $C_{\text{seg}}$ .

method	ps-mask	1-way 5-shot										2-way 5-shot										learn. params.
		classification 0/1 exact ratio (%)					segmentation mIoU (%)					classification 0/1 exact ratio (%)					segmentation mIoU (%)					
		5 <sup>0</sup>	5 <sup>1</sup>	5 <sup>2</sup>	5 <sup>3</sup>	avg.	5 <sup>0</sup>	5 <sup>1</sup>	5 <sup>2</sup>	5 <sup>3</sup>	avg.	5 <sup>0</sup>	5 <sup>1</sup>	5 <sup>2</sup>	5 <sup>3</sup>	avg.	5 <sup>0</sup>	5 <sup>1</sup>	5 <sup>2</sup>	5 <sup>3</sup>	avg.	
HSNet [12]	✗	89.3	<b>90.1</b>	66.3	<b>90.7</b>	84.1	12.5	24.7	19.4	18.1	18.7	81.3	78.4	44.0	<b>81.4</b>	71.3	13.0	25.4	22.2	18.7	19.8	2.6M
ASNet [9]	✗	84.3	89.1	66.2	90.0	82.4	11.5	22.0	14.0	17.4	16.2	72.5	<b>80.6</b>	41.8	76.8	67.9	8.7	23.1	11.8	18.0	15.4	1.3M
CST	✗	88.8	85.1	63.8	88.7	81.6	13.1	21.6	15.3	17.6	16.9	<b>88.8</b>	74.2	41.6	78.9	70.9	13.1	22.3	15.6	17.5	17.1	0.4M
DINO [2]	◇	-	-	-	-	-	20.1	23.6	16.4	16.8	19.2	-	-	-	-	-	12.9	11.9	8.4	9.4	10.7	0
CST	✓	<b>92.7</b>	89.4	<b>70.3</b>	89.2	<b>85.4</b>	<b>42.1</b>	<b>40.8</b>	<b>30.8</b>	<b>31.2</b>	<b>36.2</b>	86.2	77.4	<b>48.5</b>	73.9	<b>71.5</b>	<b>40.9</b>	<b>40.1</b>	<b>29.8</b>	<b>31.3</b>	<b>35.5</b>	0.4M

**Table 10.** Comparing 5-shot performance, *i.e.*, 5 support images per class, on FS-CS trained with image-level supervision on Pascal-5<sup>i</sup>.

method	ps-mask	1-way 5-shot										2-way 5-shot										learn. params.
		classification 0/1 exact ratio (%)					segmentation mIoU (%)					classification 0/1 exact ratio (%)					segmentation mIoU (%)					
		5 <sup>0</sup>	5 <sup>1</sup>	5 <sup>2</sup>	5 <sup>3</sup>	avg.	5 <sup>0</sup>	5 <sup>1</sup>	5 <sup>2</sup>	5 <sup>3</sup>	avg.	5 <sup>0</sup>	5 <sup>1</sup>	5 <sup>2</sup>	5 <sup>3</sup>	avg.	5 <sup>0</sup>	5 <sup>1</sup>	5 <sup>2</sup>	5 <sup>3</sup>	avg.	
CST	✗	<b>79.3</b>	<b>83.4</b>	83.2	<b>86.3</b>	83.1	11.9	11.2	8.4	11.1	10.6	<b>64.9</b>	<b>73.6</b>	72.3	72.6	70.9	10.7	11.2	8.3	10.7	10.2	0.4M
DINO [2]	◇	-	-	-	-	-	13.9	12.3	10.4	11.9	12.1	-	-	-	-	-	7.9	7.0	6.6	7.0	7.1	0
CST	✓	78.8	83.3	<b>86.7</b>	84.9	<b>83.4</b>	<b>20.8</b>	<b>20.9</b>	<b>20.6</b>	<b>21.1</b>	<b>20.8</b>	64.3	71.7	<b>77.3</b>	<b>72.8</b>	<b>71.5</b>	<b>19.0</b>	<b>21.0</b>	<b>20.8</b>	<b>20.9</b>	<b>20.4</b>	0.4M

**Table 11.** Comparing 5-shot performance, *i.e.*, 5 support images per class, on FS-CS trained with image-level supervision on COCO-20<sup>i</sup>.

method	ps-mask	1-way 1-shot										2-way 1-shot										learn. params.
		classification 0/1 exact ratio (%)					segmentation mIoU (%)					classification 0/1 exact ratio (%)					segmentation mIoU (%)					
		5 <sup>0</sup>	5 <sup>1</sup>	5 <sup>2</sup>	5 <sup>3</sup>	avg.	5 <sup>0</sup>	5 <sup>1</sup>	5 <sup>2</sup>	5 <sup>3</sup>	avg.	5 <sup>0</sup>	5 <sup>1</sup>	5 <sup>2</sup>	5 <sup>3</sup>	avg.	5 <sup>0</sup>	5 <sup>1</sup>	5 <sup>2</sup>	5 <sup>3</sup>	avg.	
CST	✗	<b>74.2</b>	<b>78.4</b>	65.9	<b>79.6</b>	74.5	11.8	10.7	8.1	10.6	10.3	<b>60.4</b>	60.2	67.5	<b>61.2</b>	62.3	10.3	10.4	7.8	9.4	9.5	0.4M
DINO [2]	◇	-	-	-	-	-	13.9	12.2	10.4	11.9	12.1	-	-	-	-	-	7.8	7.4	6.8	7.5	7.4	0
CST	✓	74.0	<b>78.4</b>	<b>82.1</b>	78.1	<b>78.2</b>	<b>20.2</b>	<b>19.8</b>	<b>19.1</b>	<b>19.5</b>	<b>19.6</b>	59.8	<b>60.5</b>	<b>68.2</b>	<b>61.2</b>	<b>62.4</b>	<b>19.0</b>	<b>17.5</b>	<b>18.1</b>	<b>18.5</b>	<b>18.3</b>	0.4M

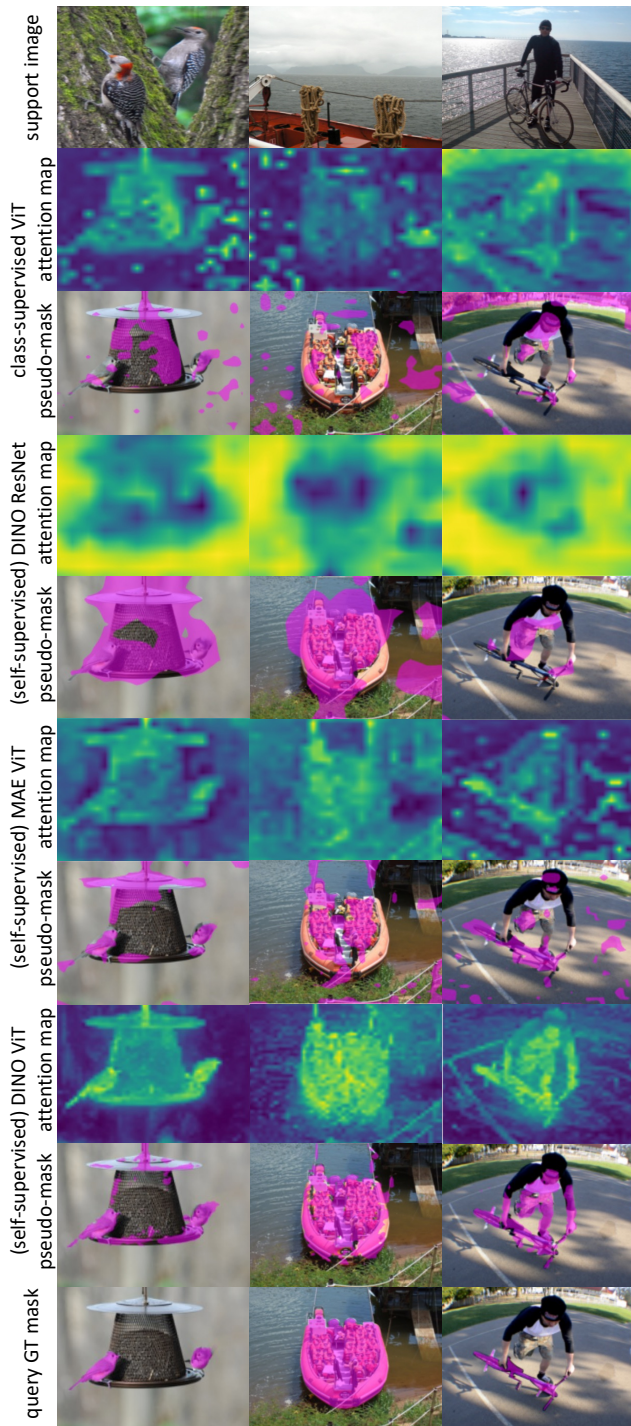
**Table 12.** Four-fold results on FS-CS with image-level supervision on COCO-20<sup>i</sup>. The results correspond to Table 2 in the main paper.

method	1-way 1-shot										2-way 1-shot									
	classification 0/1 exact ratio (%)					segmentation mIoU (%)					classification 0/1 exact ratio (%)					segmentation mIoU (%)				
	20 <sup>0</sup>	20 <sup>1</sup>	20 <sup>2</sup>	20 <sup>3</sup>	avg.	20 <sup>0</sup>	20 <sup>1</sup>	20 <sup>2</sup>	20 <sup>3</sup>	avg.	20 <sup>0</sup>	20 <sup>1</sup>	20 <sup>2</sup>	20 <sup>3</sup>	avg.	20 <sup>0</sup>	20 <sup>1</sup>	20 <sup>2</sup>	20 <sup>3</sup>	avg.
PANet [20]	64.3	66.5	68.0	67.9	66.7	25.5	24.7	25.7	24.7	25.2	42.5	49.9	53.6	47.8	48.5	24.9	25.0	23.3	21.4	23.6
PFENet [18]	70.7	70.6	71.2	72.9	71.4	30.6	34.8	29.4	32.6	31.9	35.6	34.3	43.1	32.8	36.5	23.3	23.8	20.2	23.1	22.6
HSNet [12]	74.7	77.2	78.5	77.6	77.0	36.2	34.3	32.9	34.0	34.3	57.7	62.4	67.1	<b>62.6</b>	62.5	28.9	29.6	30.3	29.3	29.5
ASNet [9]	76.2	78.8	79.2	80.2	78.6	35.7	36.8	35.3	35.6	35.8	59.5	61.5	<b>68.8</b>	62.4	63.1	29.8	33.0	33.4	30.4	31.6
CST	<b>77.6</b>	<b>82.0</b>	<b>83.1</b>	<b>80.5</b>	<b>80.8</b>	<b>36.3</b>	<b>38.3</b>	<b>37.8</b>	<b>40.7</b>	<b>38.3</b>	<b>61.0</b>	<b>66.0</b>	68.2	60.5	<b>64.0</b>	<b>34.7</b>	<b>37.1</b>	<b>36.8</b>	<b>36.3</b>	<b>36.2</b>

**Table 13.** Four-fold results on FS-CS with pixel-level supervision on COCO-20<sup>i</sup>. The results correspond to Table 5 in the main paper.

method	1-way 1-shot							1-way 5-shot							# learn. params.
	20 <sup>0</sup>	20 <sup>1</sup>	20 <sup>2</sup>	20 <sup>3</sup>	mIoU	FBIoU		20 <sup>0</sup>	20 <sup>1</sup>	20 <sup>2</sup>	20 <sup>3</sup>	mIoU	FBIoU		
RPMM [25]	29.5	36.8	28.9	27.0	30.6	-	33.8	42.0	33.0	33.3	35.5	-	38.6	M	
RePRI [1]	31.2	38.1	33.3	33.0	34.0	-	38.5	46.2	40.0	43.6	42.1	-	-		
SSP [6]	35.5	39.6	37.9	36.7	37.4	-	40.6	47.0	45.1	43.9	44.1	-	8.7M		
MMNet [22]	34.9	41.0	37.2	37.0	37.5	-	37.0	40.3	39.3	36.0	38.2	-	10.4	M	
MLC [26]	<b>46.8</b>	35.3	26.2	27.1	33.9	-	<b>54.1</b>	41.2	34.1	33.1	40.6	-	8.7	M	
NTRENet [11]	36.8	42.6	39.9	37.9	39.3	68.5	38.2	44.1	40.4	38.4	40.3	69.2	-		
CMN [23]	37.9	44.8	38.7	35.6	39.3	61.7	42.0	50.5	41.0	38.9	43.1	63.3	-		
HSNet [12]	36.3	43.1	38.7	38.7	39.2	68.2	43.3	51.3	48.2	45.0	46.9	70.7	2.6	M	
DACM [24]	37.5	44.3	40.6	40.1	40.6	68.9	44.6	<b>52.0</b>	49.2	46.4	48.1	71.6	-		
ASNet [9]	41.5	44.1	42.8	40.6	42.2	68.8	47.6	50.1	47.7	46.4	47.9	71.6	1.3	M	
CST	39.6	<b>45.8</b>	<b>45.0</b>	<b>45.5</b>	<b>44.0</b>	<b>70.3</b>	42.8	51.6	<b>50.2</b>	<b>50.2</b>	<b>48.7</b>	<b>73.7</b>	0.4	M	

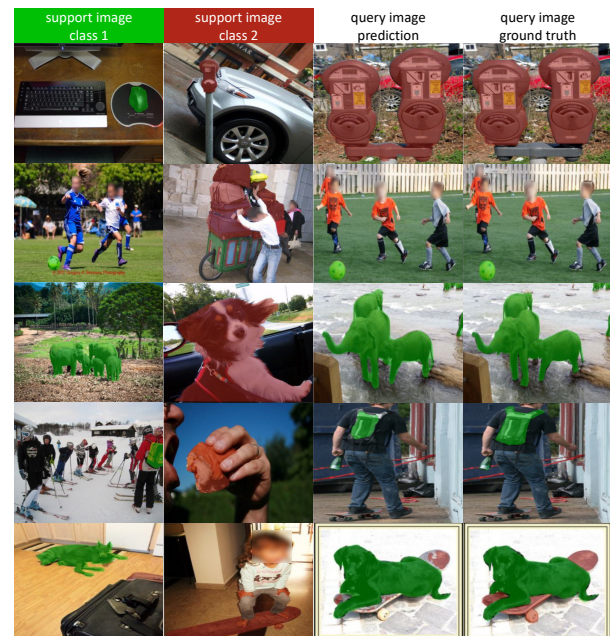
**Table 14.** Four-fold performance on the conventional few-shot segmentation task (FS-S) with image-level supervision on COCO-20<sup>i</sup> [14]. The results correspond to Table 8 in the main paper.



**Figure 3.** Pseudo-GT masks generated from class-supervised ViT [3], self-supervised DINO ResNet [2,8], self-supervised MAE ViT [3, 7], and self-supervised DINO ViT [2, 3] from the top.



**Figure 4.** 2-way 1-shot segmentation prediction of CST trained with image-level labels. Image frames on the support images distinguish classes by colors.



**Figure 5.** 2-way 1-shot segmentation prediction of CST. The model is trained with pixel-level labels during training, and pixel-level support annotations are also given during testing (as overlaid on the support images with colors). Human faces are anonymized for visualization.

## 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 *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 4
- [2] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2021. 1, 2, 4, 5
- [3] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *Proc. International Conference on Learning Representations (ICLR)*, 2021. 2, 5
- [4] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International Journal of Computer Vision (IJCV)*, 2010. 1
- [5] William Falcon et al. Pytorch lightning. *GitHub. Note: <https://github.com/PyTorchLightning/pytorch-lightning>*, 2019. 1
- [6] Qi Fan, Wenjie Pei, Yu-Wing Tai, and Chi-Keung Tang. Self-support few-shot semantic segmentation. In *Proc. European Conference on Computer Vision (ECCV)*, 2022. 4
- [7] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 5
- [8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2, 5
- [9] Dahyun Kang and Minsu Cho. Integrative few-shot learning for classification and segmentation. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 3, 4
- [10] 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 *Proc. European Conference on Computer Vision (ECCV)*, 2014. 1
- [11] Yuanwei Liu, Nian Liu, Qinglong Cao, Xiwen Yao, Junwei Han, and Ling Shao. Learning non-target knowledge for few-shot semantic segmentation. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 4
- [12] Juhong Min, Dahyun Kang, and Minsu Cho. Hypercorrelation squeeze for few-shot segmentation. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2021. 2, 3, 4
- [13] Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In *International Conference on Machine Learning (ICML)*, 2010. 1
- [14] Khoi Nguyen and Sinisa Todorovic. Feature weighting and boosting for few-shot segmentation. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2019. 4
- [15] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning internal representations by error propagation. Technical report, California Univ San Diego La Jolla Inst for Cognitive Science, 1985. 2
- [16] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. 1
- [17] Amirreza Shaban, Shray Bansal, Zhen Liu, Irfan Essa, and Byron Boots. One-shot learning for semantic segmentation. In *Proc. British Machine Vision Conference (BMVC)*, 2017. 2
- [18] Zhuotao Tian, Hengshuang Zhao, Michelle Shu, Zhicheng Yang, Ruiyu Li, and Jiaya Jia. Prior guided feature enrichment network for few-shot segmentation. In *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2020. 2, 3, 4
- [19] Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and composing robust features with denoising autoencoders. In *Proc. International Conference on Machine Learning (ICML)*, 2008. 2
- [20] Kaixin Wang, Jun Hao Liew, Yingtian Zou, Daquan Zhou, and Jiashi Feng. Panet: Few-shot image semantic segmentation with prototype alignment. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2019. 3, 4
- [21] Yuxin Wu and Kaiming He. Group normalization. In *Proc. European Conference on Computer Vision (ECCV)*, 2018. 1
- [22] Zhonghua Wu, Xiangxi Shi, Guosheng Lin, and Jianfei Cai. Learning meta-class memory for few-shot semantic segmentation. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2021. 4
- [23] Guo-Sen Xie, Huan Xiong, Jie Liu, Yazhou Yao, and Ling Shao. Few-shot semantic segmentation with cyclic memory network. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2021. 4
- [24] Zhitong Xiong, Haopeng Li, and Xiao Xiang Zhu. Doubly deformable aggregation of covariance matrices for few-shot segmentation. In *Proc. European Conference on Computer Vision (ECCV)*, 2022. 4
- [25] Boyu Yang, Chang Liu, Bohao Li, Jianbin Jiao, and Ye Qixiang. Prototype mixture models for few-shot semantic segmentation. In *Proc. European Conference on Computer Vision (ECCV)*, 2020. 4
- [26] Lihe Yang, Wei Zhuo, Lei Qi, Yinghuan Shi, and Yang Gao. Mining latent classes for few-shot segmentation. In *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2021. 4