

# 01 Few-Shot Segmentation via Cycle-Consistent Transformer

## Few-Shot Segmentation via Cycle-Consistent Transformer

Gengwei Zhang<sup>1,2\*</sup>, Guoliang Kang<sup>3</sup>, Yi Yang<sup>4</sup>, Yunchao Wei<sup>5,6†</sup>

<sup>1</sup> Baidu Research

<sup>2</sup> ReLER, Centre for Artificial Intelligence, University of Technology Sydney

<sup>3</sup> University of Texas, Austin

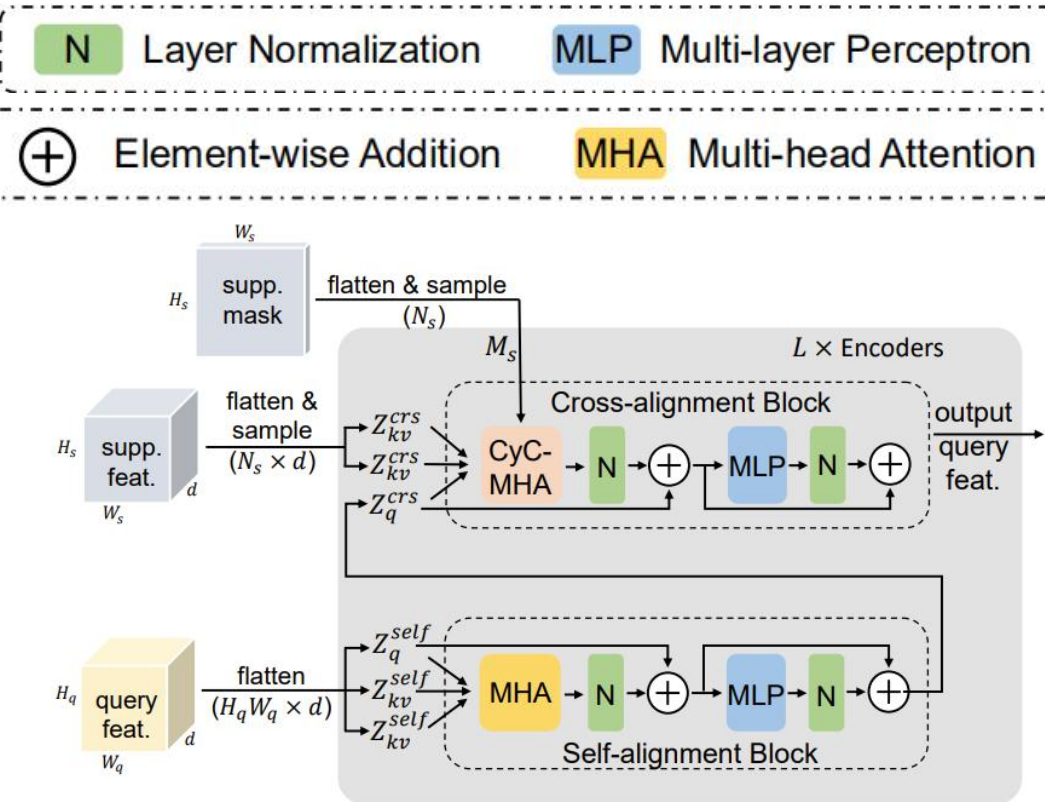
<sup>4</sup> CCAI, College of Computer Science and Technology, Zhejiang University

<sup>5</sup> Institute of Information Science, Beijing Jiaotong University

<sup>6</sup> Beijing Key Laboratory of Advanced Information Science and Network

{zgw david, kgl.prml, wychao1987, yee.i.yang}@gmail.com

- Neural Information Processing Systems (NIPS 2021)

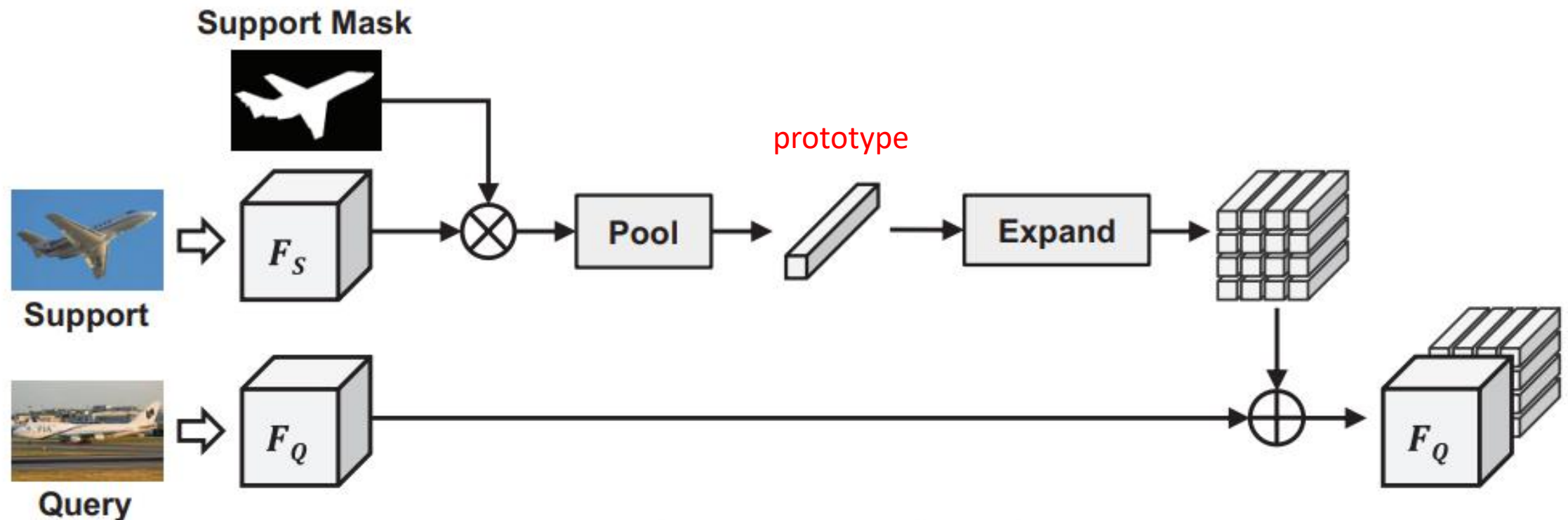


# 01 Background for Few-Shot Segmentation

## • Motivation

- Deep learning based computer vision systems have largely depended on large-scale training sets
- Deep networks mostly work with predefined classes and are incapable of generalizing to new ones

➔ **Few shot:** learn how to recognize novel objects after seeing only a handful of exemplars



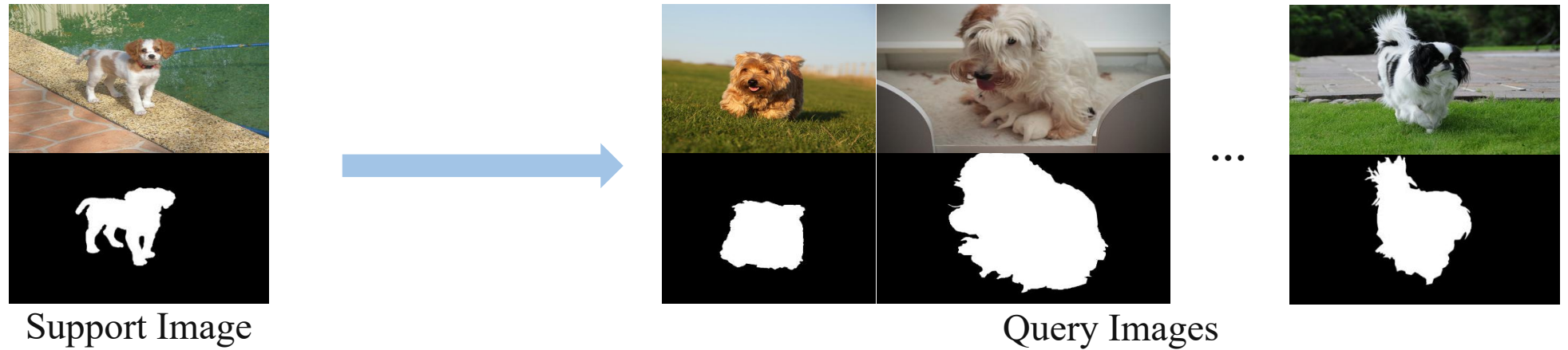
Feature from Res2 + Res3

# 01 Background for Few-Shot Segmentation

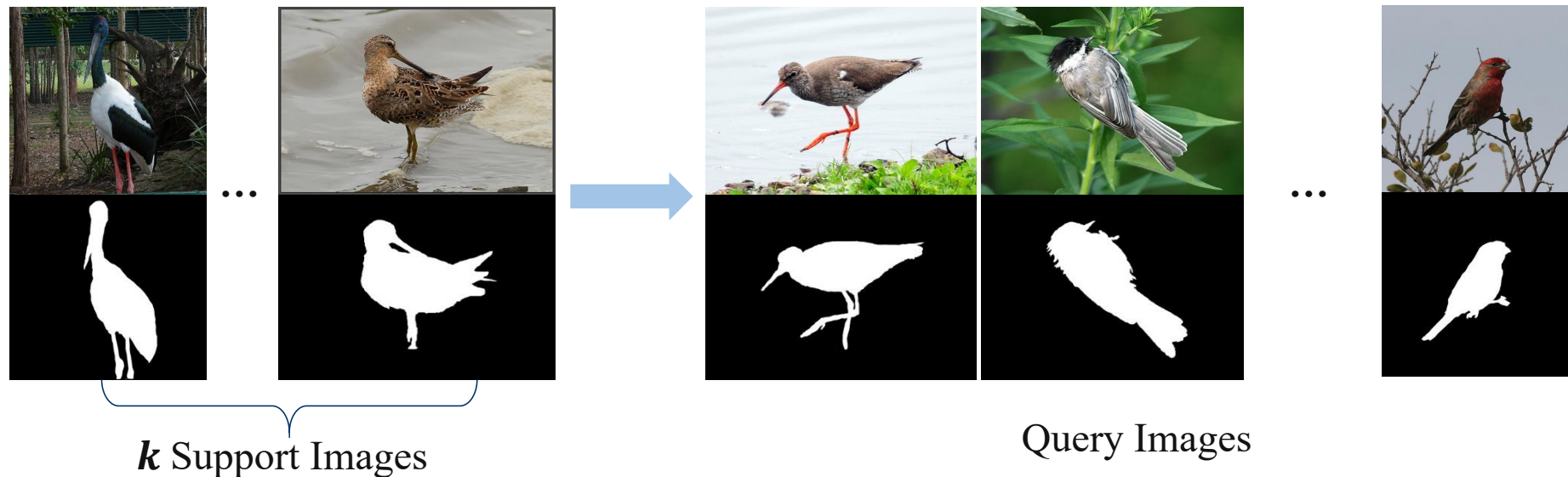
- **Implementation**

- Divide the dataset into  $\{?, ?, ?, ?, ?\}$  and  $\{?, ?, ?, ?, ?, ?, ?\}$ ,  $\{?, ?, ?, ?, ?\} \cap \{?, ?, ?, ?, ?, ?\} = \emptyset$
- **One-shot segmentation** and  **$k$ -shot segmentation**

One-shot Segmentation:



$k$ -shot Segmentation:





# 01 Few-Shot Segmentation via Cycle-Consistent Transformer

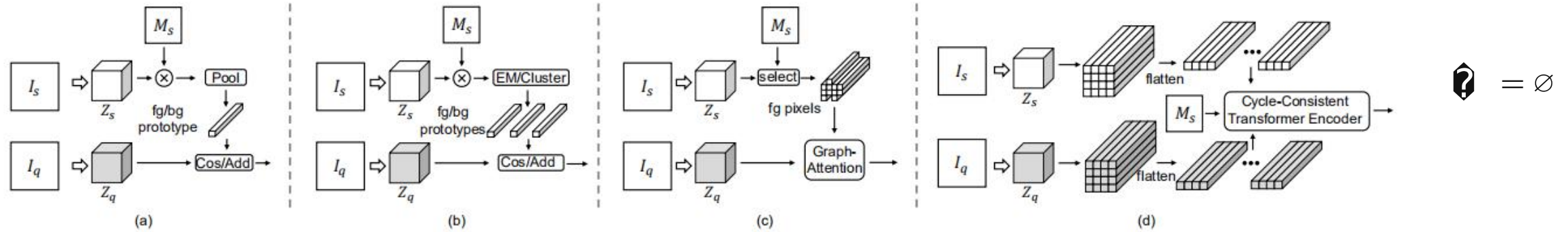
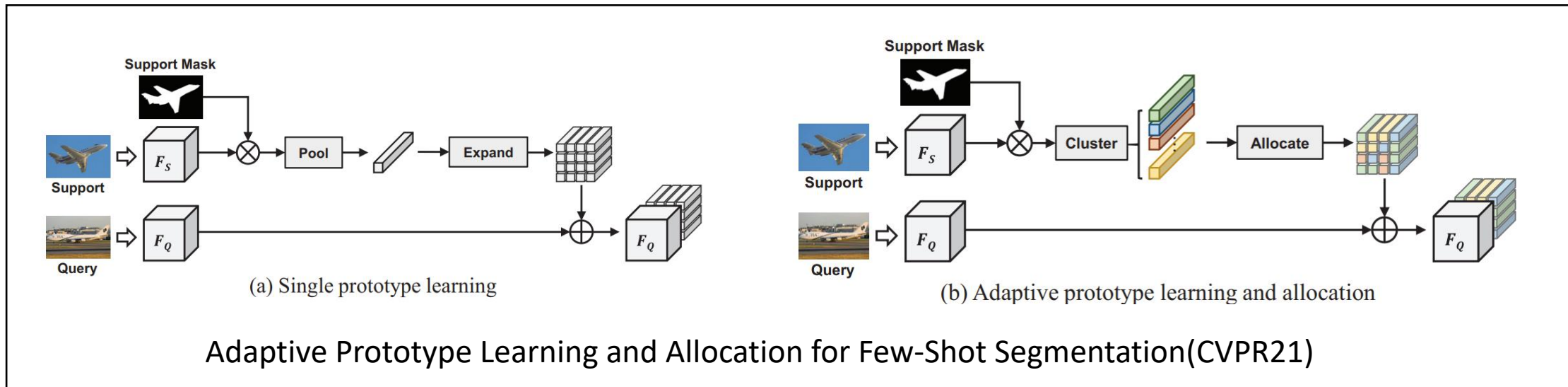


Figure 1: Different learning frameworks for few-shot segmentation, from the perspective of ways to utilize support information. (a) Class-wise mean pooling based method. (b) Clustering based method. (c) Foreground pixel attention method. (d) Our Cycle-Consistent TRansformer (CyCTR) framework that enables all beneficial support pixel-level features (foreground and background) to be considered.



# 01 Few-Shot Segmentation via Cycle-Consistent Transformer

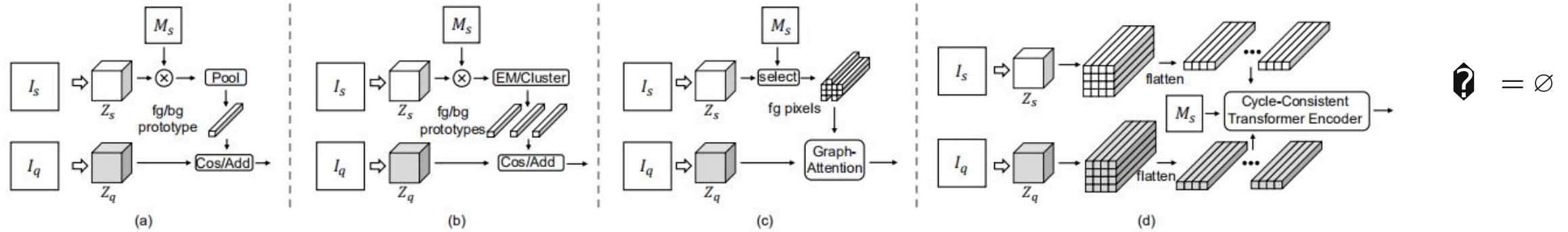
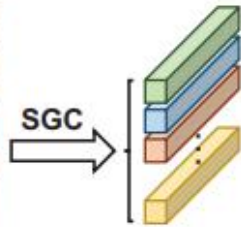
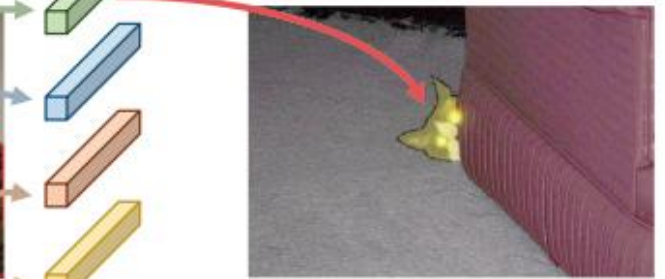
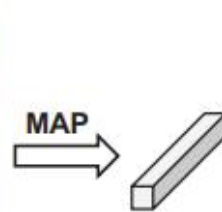


Figure 1: Different learning frameworks for few-shot segmentation, from the perspective of ways to utilize support information. (a) Class-wise mean pooling based method. (b) Clustering based method. (c) Foreground pixel attention method. (d) Our Cycle-Consistent TRansformer (CyCTR) framework that enables all beneficial support pixel-level features (foreground and background) to be considered.



(a) SGC is adaptive to object scale variation



(b) GPA is adaptive to object shape variation

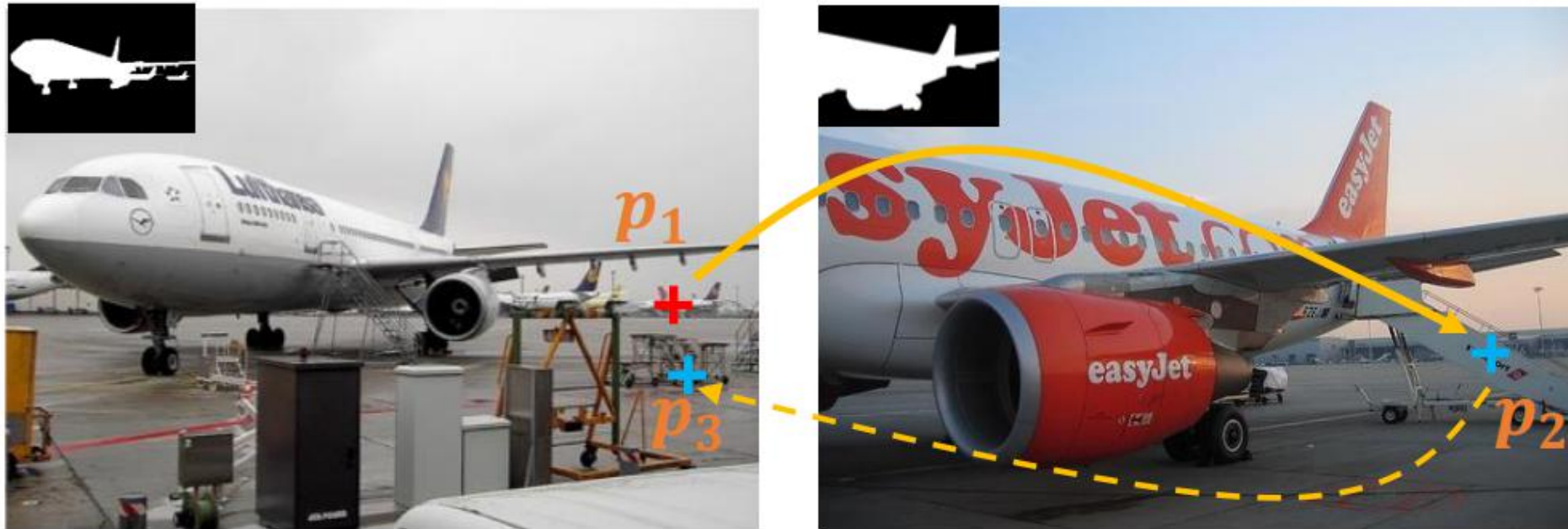
Adaptive Prototype Learning and Allocation for Few-Shot Segmentation(CVPR21)



# 01 Few-Shot Segmentation via Cycle-Consistent Transformer

- **Motivation**

- Many pixel-level **support features** are quite different from the **query ones**, and thus may confuse the attention.



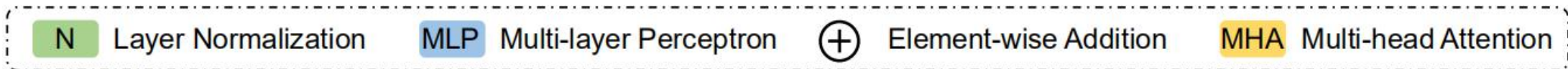
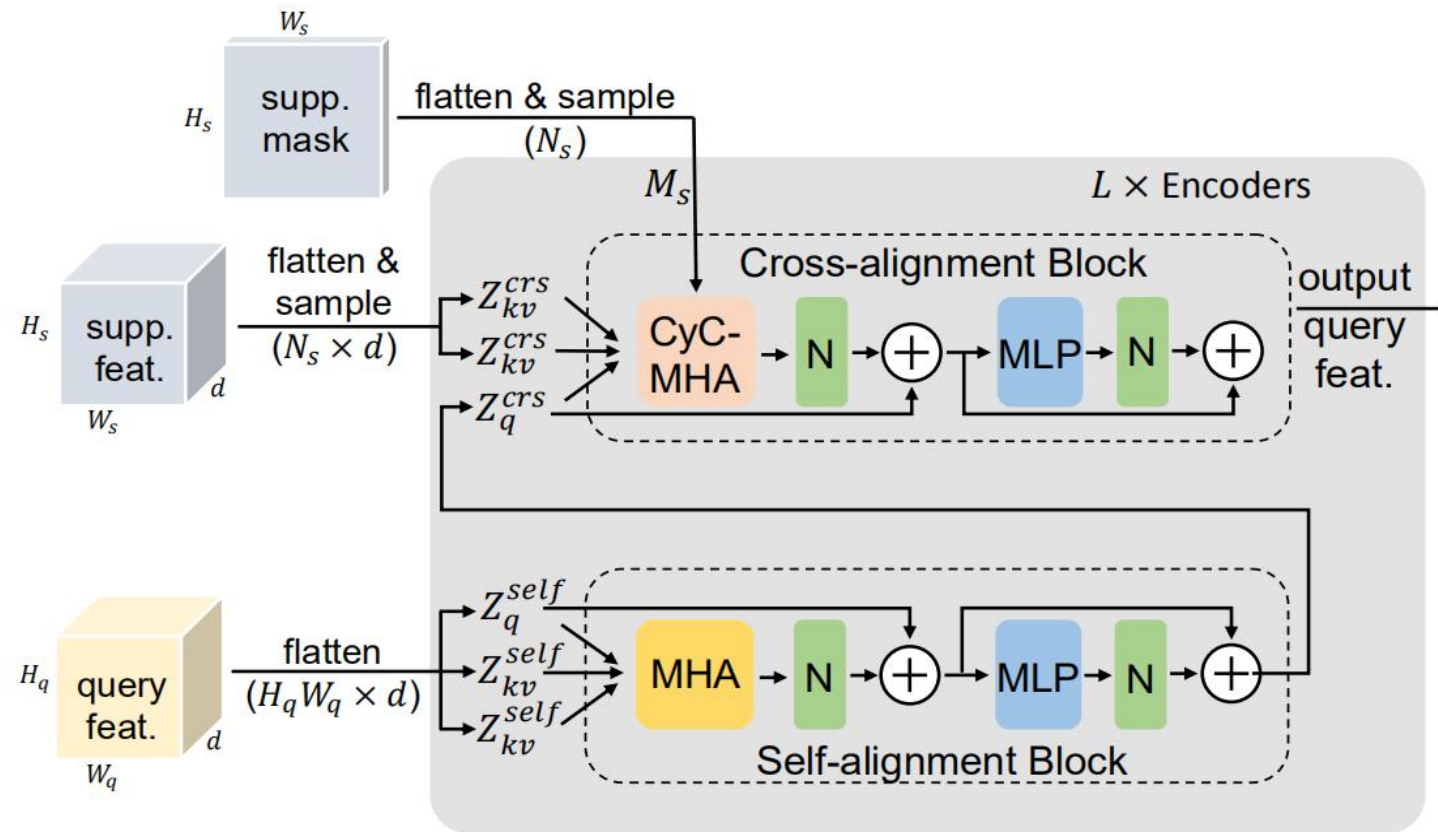
support image

query image



# 01 Few-Shot Segmentation via Cycle-Consistent Transformer

- Framework of the proposed Cycle-Consistent Transformer (CyCTR)
  - Self-alignment block for utilizing global context within the query feature map
  - Cross-alignment block for aggregate information from support images



# 01 Few-Shot Segmentation via Cycle-Consistent Transformer

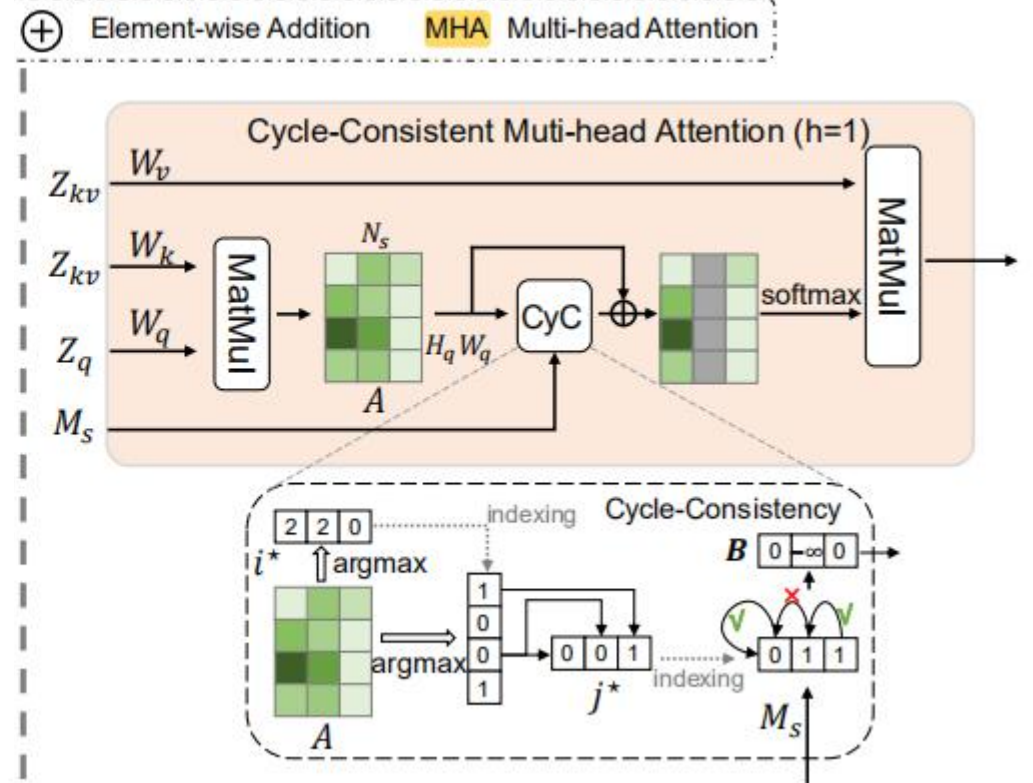
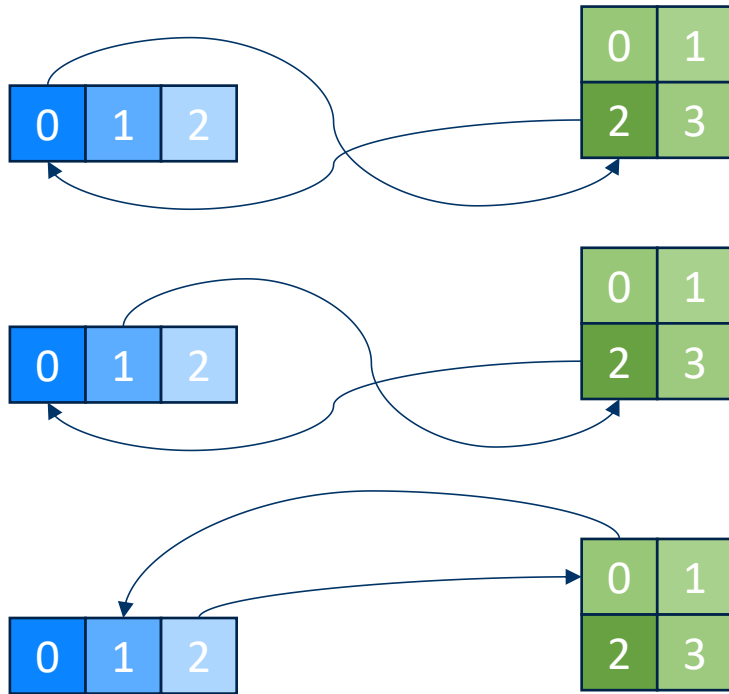
- Cycle-Consistent Attention ( Cross-alignment block )

$$A = \frac{QK^T}{\sqrt{d}}, A \in \mathbb{R}^{H_q W_q \times N_s}$$

$$i^* = \operatorname{argmax}_i A_{(i,j)}, \quad B_j = \begin{cases} 0, & \text{if } M_s(j) = M_s(j^*) \\ -\infty, & \text{if } M_s(j) \neq M_s(j^*) \end{cases}$$

$$j^* = \operatorname{argmax}_j A_{(i^*,j)}$$

$$\text{CyCAtten}(Q_i, K_i, V_i) = \text{softmax}(A_{(i)} + B)V,$$





# 01 Few-Shot Segmentation via Cycle-Consistent Transformer

- Mask-guided sparse sampling and K-shot Setting

$$N_{fg} \leq \frac{N_s}{2}, N_s \leq kH_sW_s$$

$$N_{bg} = N_s - N_{fg}$$

With a proper  $N_s$ , the sampling operation reduces the computational complexity, and makes our algorithm more scalable with the increase of spatial size of support images.

# 01 Few-Shot Segmentation via Cycle-Consistent Transformer

- **Self-alignment block**
  - Referring to deformable detr

$$\text{PredAtten}(Q_r, V_r) = \sum_g^P \text{softmax}(A')_{(r,g)} V_{r+\Delta_{(r,g)}},$$

$$\text{Atten}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V,$$

Compare with original attention

$$\Delta = f(Q + \text{Coord}) \text{ and } A' = g(Q + \text{Coord}), \quad A' \in \mathbb{R}^{H_q W_q \times P}$$

$f(\cdot)$  and  $g(\cdot)$  are two fully connected layers that predict the offsets<sup>3</sup> and attention weights.

# 01 Few-Shot Segmentation Setting

- **Dataset**

- **Pascal-5<sup>i</sup>**: 20 categories, 15 classes are used for training and 5 classes for test.
- **COCO- 20<sup>i</sup>**: 80 categories, 60 classes are used for training and 20 classes for test.

- **Evaluation Metric**

- mIoU
- FB-IoU



# 01 Few-Shot Segmentation via Cycle-Consistent Transformer

- Results

Table 1: Comparison with other state-of-the-art methods for 1-shot and 5-shot segmentation on PASCAL-5<sup>i</sup> using the mIoU (%) evaluation metric. Best results are shown in bold.

Method	Backbone	1-shot					5-shot				
		5 <sup>0</sup>	5 <sup>1</sup>	5 <sup>2</sup>	5 <sup>3</sup>	Mean	5 <sup>0</sup>	5 <sup>1</sup>	5 <sup>2</sup>	5 <sup>3</sup>	Mean
PANet [35]	Vgg-16	42.3	58.0	51.1	41.2	48.1	51.8	64.6	59.8	46.5	55.7
FWB [23]		47.0	59.6	52.6	48.3	51.9	50.9	62.9	56.5	50.1	55.1
SG-One [45]		40.2	58.4	48.4	38.4	46.3	41.9	58.6	48.6	39.4	47.1
RPMM [41]		47.1	65.8	50.6	48.5	53.0	50.0	66.5	51.9	47.6	54.0
CANet [44]	Res-50	52.5	65.9	51.3	51.9	55.4	55.5	67.8	51.9	53.2	57.1
PGNet [43]		56.0	66.9	50.6	50.4	56.0	57.7	68.7	52.9	54.6	58.5
RPMM [41]		55.2	66.9	52.6	50.7	56.3	56.3	67.3	54.5	51.0	57.3
PPNet [18]		47.8	58.8	53.8	45.6	51.5	58.4	67.8	<b>64.9</b>	56.7	62.0
PFENet [30]		61.7	69.5	55.4	56.3	60.8	63.1	70.7	55.8	<b>57.9</b>	61.9
CyCTR (Ours)	Res-50	<b>67.8</b>	<b>72.8</b>	<b>58.0</b>	<b>58.0</b>	<b>64.2</b>	<b>71.1</b>	<b>73.2</b>	60.5	57.5	<b>65.6</b>
FWB [23]	Res-101	51.3	64.5	56.7	52.2	56.2	54.9	67.4	<b>62.2</b>	55.3	59.9
DAN [34]		54.7	68.6	<b>57.8</b>	51.6	58.2	57.9	69.0	60.1	54.9	60.5
PFENet [30]		60.5	69.4	54.4	55.9	60.1	62.8	70.4	54.9	57.6	61.4
CyCTR (Ours)	Res-101	<b>69.3</b>	<b>72.7</b>	56.5	<b>58.6</b>	<b>64.3</b>	<b>73.5</b>	<b>74.0</b>	58.6	<b>60.2</b>	<b>66.6</b>

# 01 Few-Shot Segmentation via Cycle-Consistent Transformer

- Results

Table 2: Comparison with other state-of-the-art methods for 1-shot and 5-shot segmentation on COCO-20<sup>i</sup> using the mIoU (%) evaluation metric. Best results are shown in bold.

Method	Backbone	1-shot					5-shot				
		20 <sup>0</sup>	20 <sup>1</sup>	20 <sup>2</sup>	20 <sup>3</sup>	Mean	20 <sup>0</sup>	20 <sup>1</sup>	20 <sup>2</sup>	20 <sup>3</sup>	Mean
FWB [23]	Res-101	19.9	18.0	21.0	28.9	21.2	19.1	21.5	23.9	30.1	23.7
PPNet [18]	Res-50	28.1	30.8	29.5	27.7	29.0	39.0	40.8	37.1	37.3	38.5
RPMM [41]	Res-50	29.5	36.8	29.0	27.0	30.6	33.8	42.0	33.0	33.3	35.5
PFENet [30]	Res-101	34.3	33.0	32.3	30.1	32.4	38.5	38.6	38.2	34.3	37.4
CyCTR (Ours)	Res-50	<b>38.9</b>	<b>43.0</b>	<b>39.6</b>	<b>39.8</b>	<b>40.3</b>	<b>41.1</b>	<b>48.9</b>	<b>45.2</b>	<b>47.0</b>	<b>45.6</b>

# 01 Few-Shot Segmentation via Cycle-Consistent Transformer

- Results

Table 4: Ablation studies that validate the effectiveness of each component in our Cycle-Consistent Transformer. The first result is obtained by our baseline (see Section 4.2 for details).

self-alignment	cross-alignment	CyCTR (pred)	CyCTR (fg. only)	CyCTR	mIoU (%)
					58.8
✓					61.6
✓	✓				61.2
✓	✓	✓			61.9
✓	✓		✓		62.0
✓	✓			✓	62.8