

## **Fine-Grained Entity Segmentation**

Lu Qi<sup>1\*</sup>, Jason Kuen<sup>2\*</sup>, Weidong Guo<sup>3\*</sup>, Tiancheng Shen<sup>4</sup>, Jiuxiang Gu<sup>2</sup>, Wenbo Li<sup>4</sup>,  
Jiaya Jia<sup>4</sup>, Zhe Lin<sup>2</sup>, Ming-Hsuan Yang<sup>1</sup>

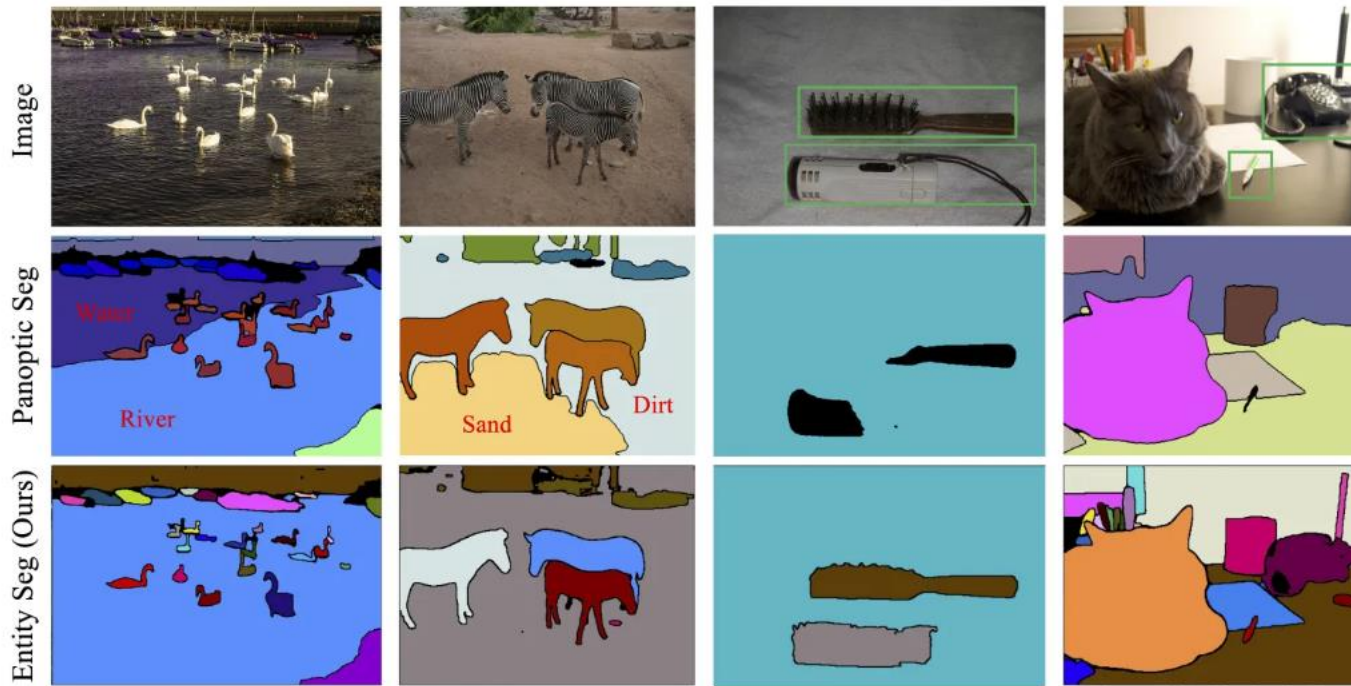
<sup>1</sup>The University of California, Merced

<sup>2</sup>Adobe Research

<sup>3</sup>QQ Browser Lab, Tencent

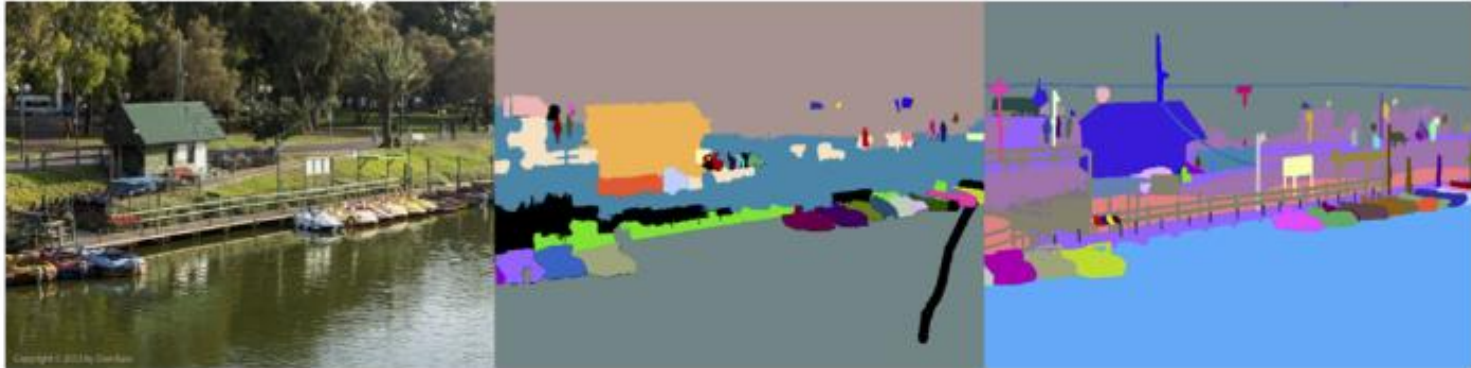
<sup>4</sup>The Chinese University of Hong Kong

# Entity Segmentation



Each entity is a thing or stuff that does not consider category information

# EntitySeg Dataset

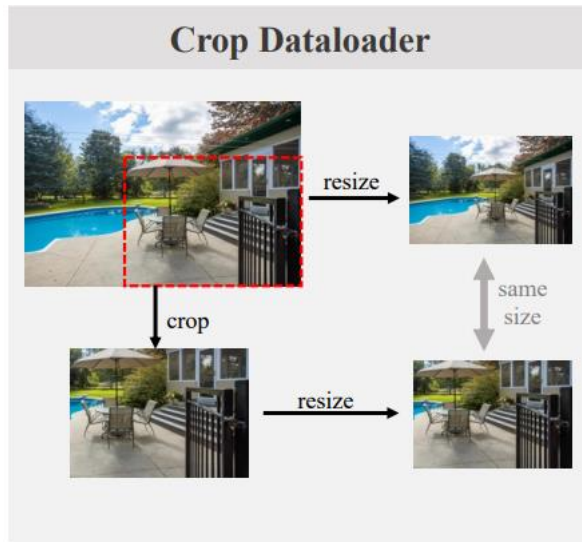


- Large-scale
- High-quality
- High-resolution



# CropFormer

Add high-resolution crop inputs to improve the quality of fine-grained entity segmentation.

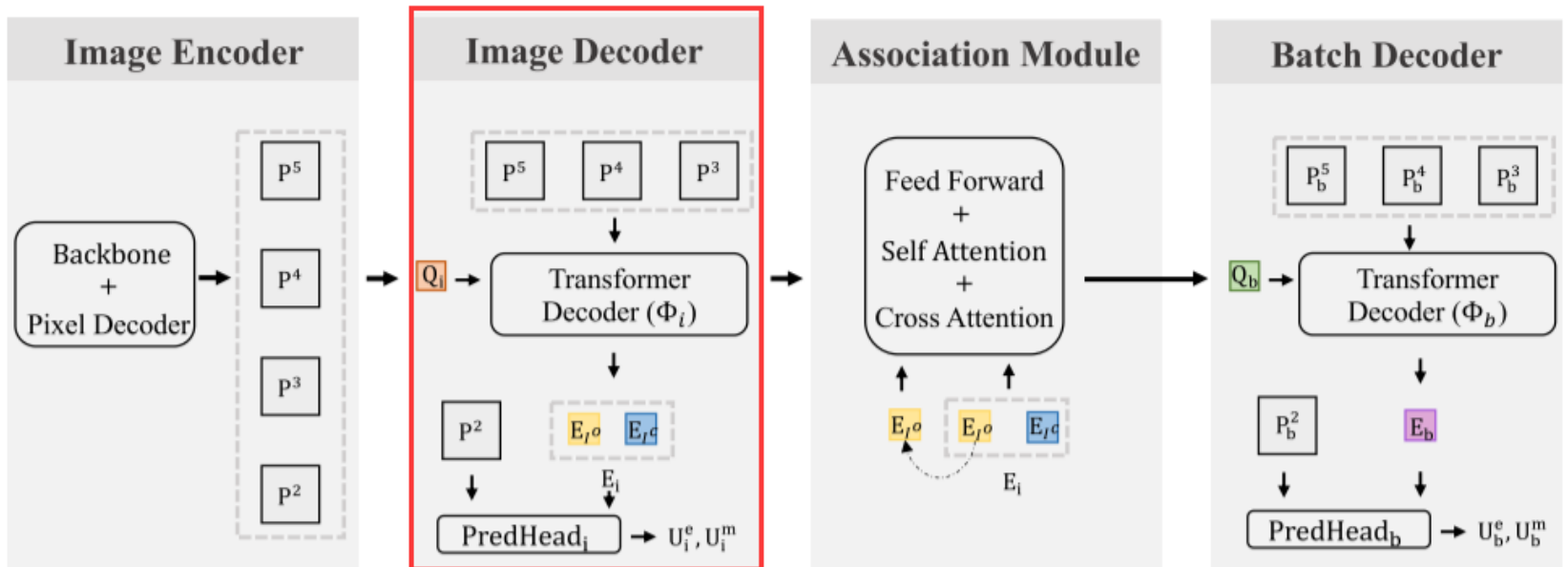


Crop & Resize



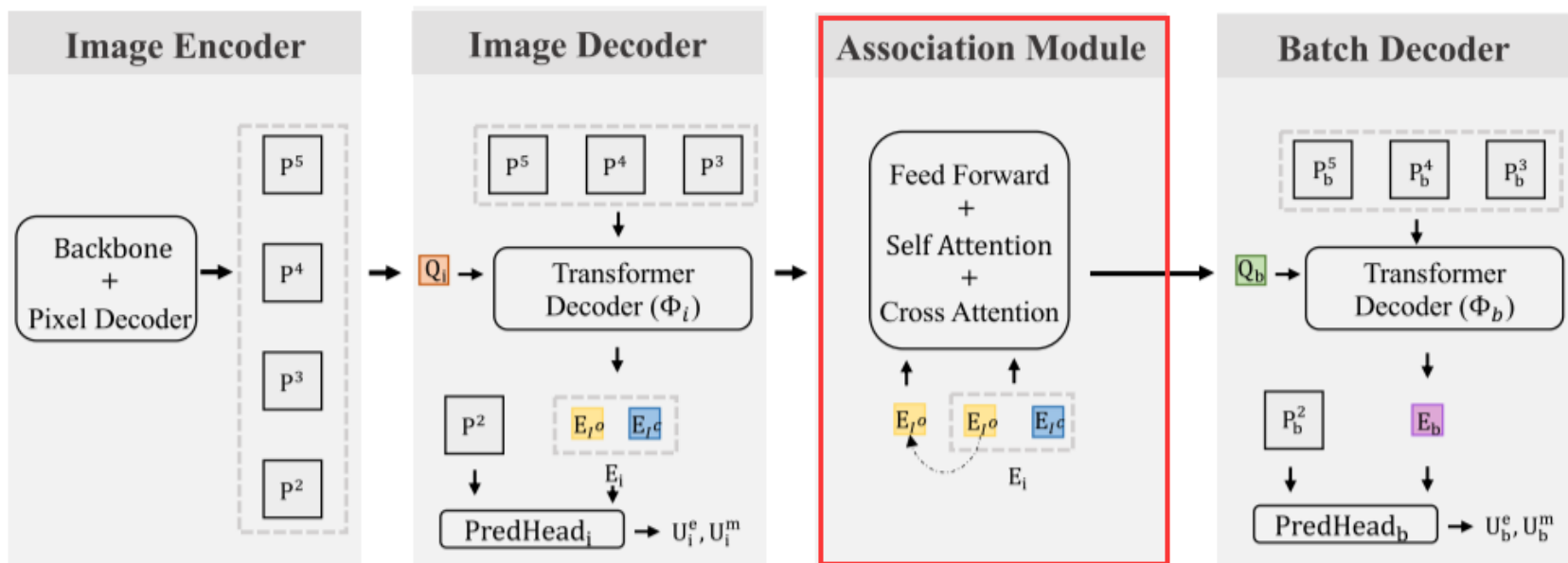
$$\mathbf{I} = \{I^O, I^C\} = \Gamma\{I^O, \delta\}$$

$$\mathbf{I} \in R^{N \times 2 \times H_I \times W_I \times 3}$$



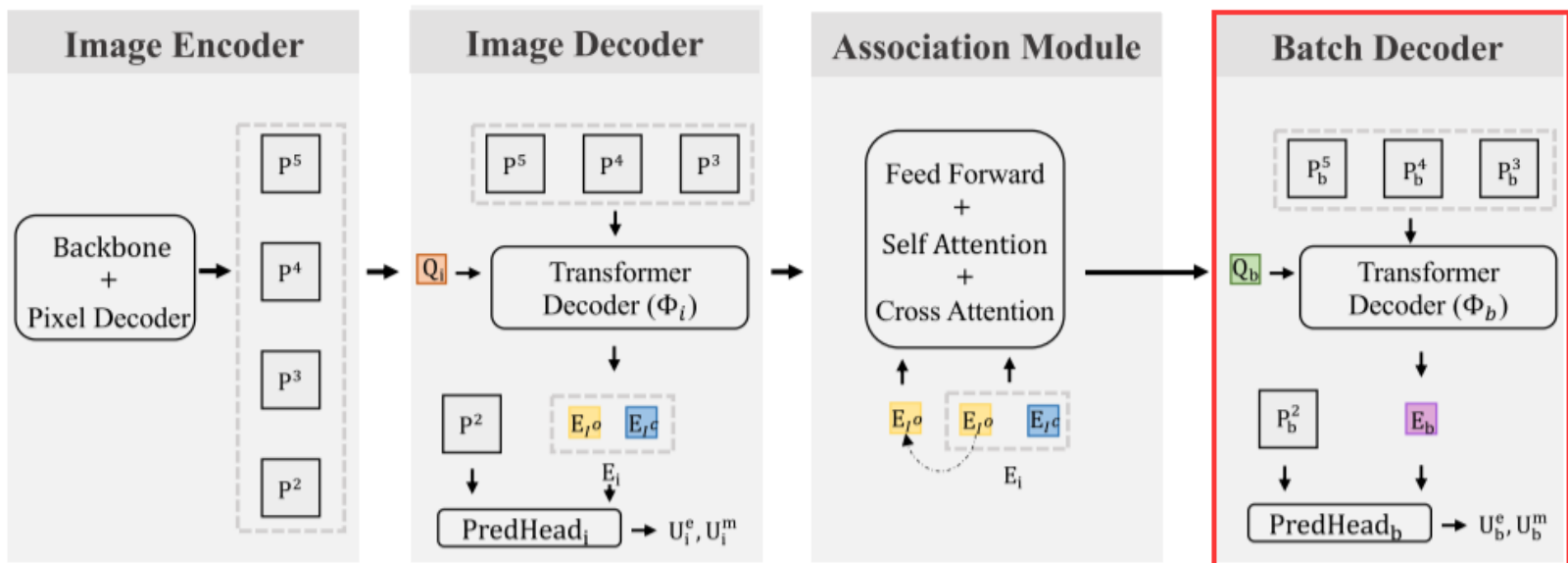
$$\mathbf{E}_i = \Phi_i(\mathbf{Q}_i, \Theta(\mathbf{I})) \quad \mathbf{E}_i \in \mathbb{R}^{N \times 2 \times 1 \times 1 \times K}$$

$$\mathbf{U}_i^e, \mathbf{U}_i^m = \text{PredHead}_i(\mathbf{E}_i, \mathbf{P}_i^h)$$



Generate batch queries  $Q_b$  that are fully shared by the full image and its crop to represent the same entities consistently

$$Q_b = \text{FFN}(\text{SAtt}(\text{XAtt}(\underbrace{f_q(\mathbf{E}_{I^o})}_{\text{query}}, \underbrace{f_k(\mathbf{E}_i)}_{\text{key}}, \underbrace{f_v(\mathbf{E}_i)}_{\text{value}})))$$



$$\mathbf{E}_b = \Phi_b(\mathbf{Q}_b, \Theta(\mathbf{I})) \quad \mathbf{Q}_b \in R^{N \times 1 \times 1 \times 1 \times K}, \quad \mathbf{E}_b \in R^{N \times 1 \times 1 \times 1 \times K}$$

↓ broadcast

$$\mathbf{U}_b^e, \mathbf{U}_b^m = \text{PredHead}_b(\mathbf{E}_b, \mathbf{P}_b^h) \quad \mathbf{E}_b \in R^{N \times 2 \times 1 \times 1 \times K}$$

$$\mathcal{L} = \sum_{\mathbf{k} \in \{\mathbf{i}, \mathbf{b}\}} \mathcal{L}_{\mathbf{k}}^{\text{ce}}(\mathbf{U}_{\mathbf{k}}^{\text{e}}, \mathbf{G}_{\mathbf{k}}^{\text{e}}) + \sum_{\mathbf{k} \in \{\mathbf{i}, \mathbf{b}\}} \mathcal{L}_{\mathbf{k}}^{\text{bce}}(\mathbf{U}_{\mathbf{k}}^{\text{m}}, \mathbf{G}_{\mathbf{k}}^{\text{m}}) + \sum_{\mathbf{k} \in \{\mathbf{i}, \mathbf{b}\}} \mathcal{L}_{\mathbf{k}}^{\text{dice}}(\mathbf{U}_{\mathbf{k}}^{\text{m}}, \mathbf{G}_{\mathbf{k}}^{\text{m}}), \quad (7)$$

two separate losses  $L_i$  and  $L_b$  for image- and batch-level predictions



Method	Decoder	AP <sup>e</sup>	AP <sub>50</sub> <sup>e</sup>	AP <sub>75</sub> <sup>e</sup>	RT (ms)
SS-Mask2Former	Image-O	39.5	56.9	40.2	637
SS-Mask2Former( $\times \delta$ )	Image-O	39.9	57.4	40.3	876
MS-Mask2Former	Image-O	39.2	56.3	39.5	1324
MS-Mask2Former	Batch-OC	39.3	56.4	39.7	2783
CropFormer	Image-O	39.3	56.7	39.8	637
	Batch-O	39.1	56.6	39.7	1514
	Batch-C	40.2	57.5	40.8	1507
	Batch-OC	<b>41.0</b>	<b>58.4</b>	<b>41.9</b>	1545

Table 7: Ablation study on the ensemble strategy on full image and four crops. The ‘Decoder’ column indicates whether we use the inference result of the full image (‘O’), four cropped patches (‘C’), or both of them (‘OC’) from the ‘Image’ or ‘Batch’ decoder. Here, the run-time (RT) is the time of network forward except the data processing and calculated on A100 GPU.

$\delta$	AP <sup>c</sup>
0.5	38.5
0.6	40.2
0.7	<b>41.0</b>
0.8	40.9

(a)

Train	Test	AP <sup>c</sup>
Random	Fixed (4)	39.7
Fixed (4)	Fixed (4)	41.0
Fixed (4)	Fixed (8)	<b>41.3</b>
Fixed (8)	Fixed (8)	41.0

(b)

XAtt	SAtt	FFN	AP <sup>c</sup>
✓	○	○	40.7
✓	○	✓	40.8
✓	✓	○	40.8
✓	✓	✓	<b>41.0</b>

(c)

Table 9: Ablation study on the usage of crop ratio  $\delta$ , crop type and association module in CropFormer. In sub-table (b), ‘Random’ indicates random crops and ‘Fixed (4/8)’ indicates 4 or 8 fixed corner crops. In sub-table(c), ✓ and ○ means whether we use the module or not.