

# PolyFormer: Referring Image Segmentation as Sequential Polygon Generation



Jiang Liu<sup>1\*†</sup>



Hui Ding<sup>2\*</sup>



Zhaowei Cai<sup>2</sup>



Yuting Zhang<sup>2</sup>



Ravi Kumar  
Satzoda<sup>2</sup>



Vijay Mahadevan<sup>2</sup>



R. Manmatha<sup>2</sup>

<sup>1</sup> Johns Hopkins University, <sup>2</sup> AWS AI Labs, \*Equal Contribution, †Work done during internship at AWS AI Labs

<https://polyformer.github.io/>



JOHNS HOPKINS  
UNIVERSITY



Code & Demo Available

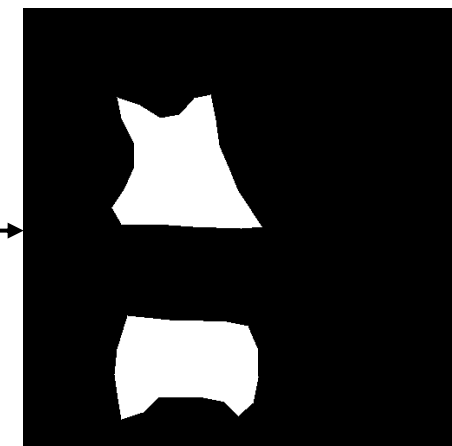
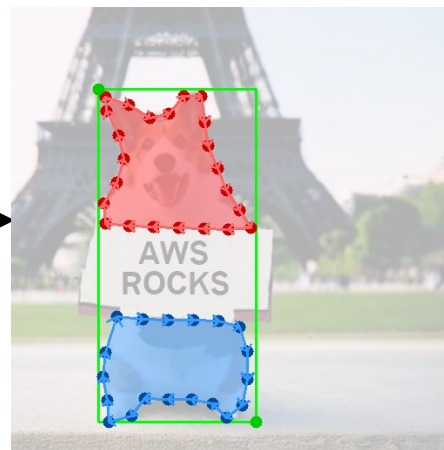
# PolyFormer Overview

- **Unified framework** for referring image segmentation and referring expression comprehension
- **Regression-based decoder** for accurate coordinate prediction
- **Superior performance** across all main referring image segmentation



*a cute corgi holding a sign that says "AWS ROCKS"*

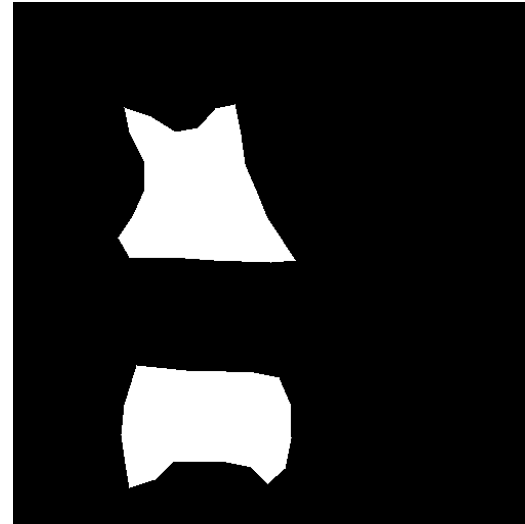

PolyFormer



# Referring Image Segmentation



RIS

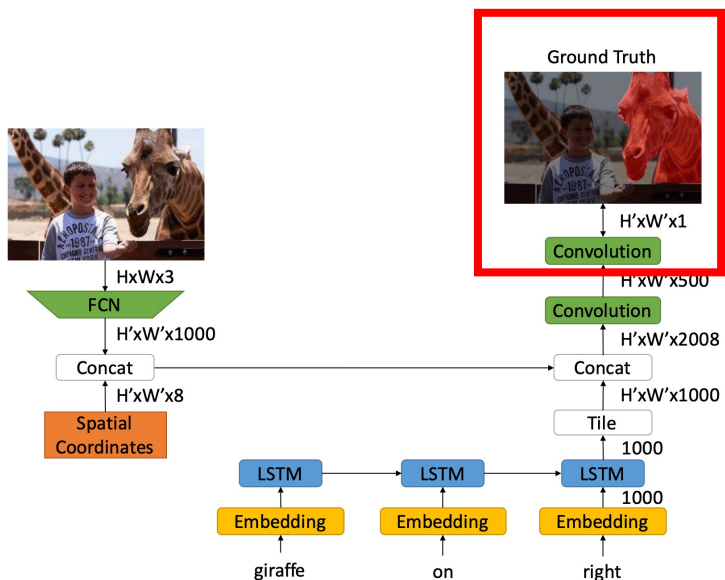


*a cute corgi holding a sign  
that says "AWS ROCKS"*

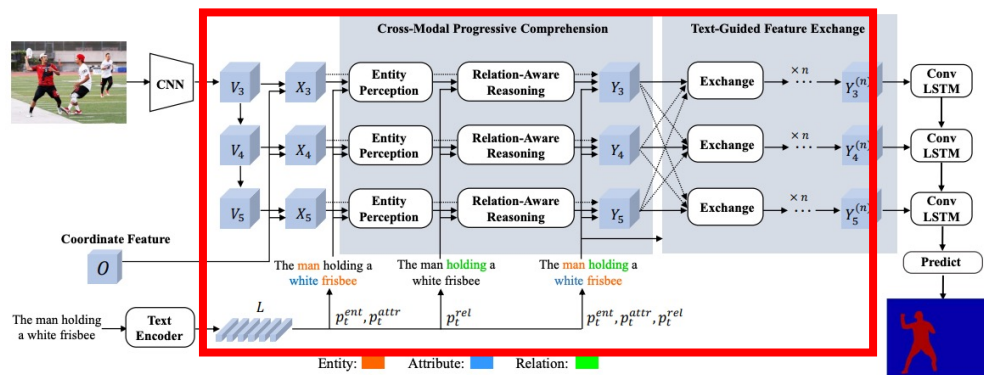


# Existing Work

- **Mask-based** dense prediction
  - Neglect the structure among the output predictions
- Complex **multi-modal feature fusion**



RMI [Liu et al, ICCV2017]



Huang et al. "Referring Image Segmentation via Cross-Modal Progressive Comprehension." CVPR2020



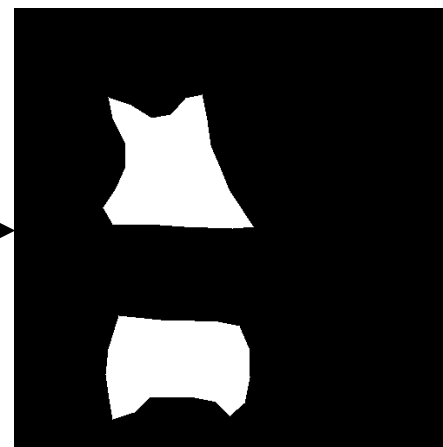
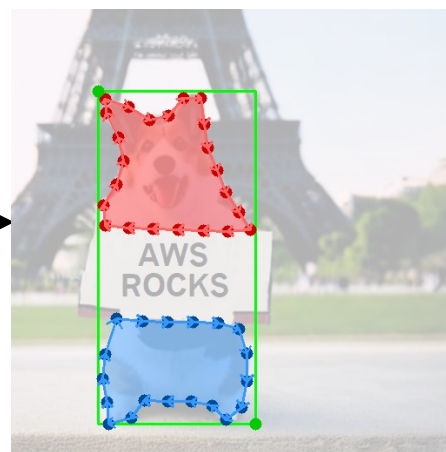
# PolyFormer

- Sequence-to-sequence formulation

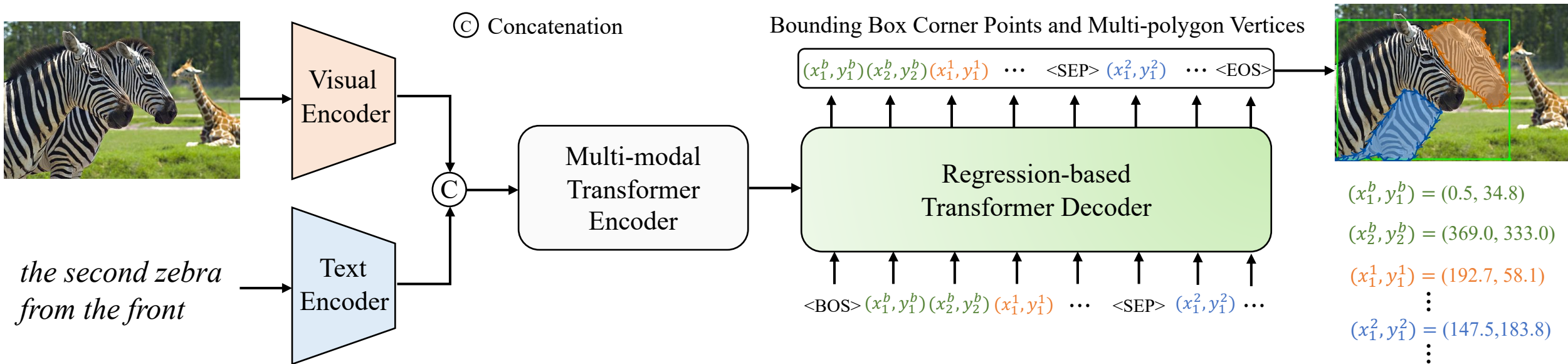


*a cute corgi holding a sign  
that says "AWS ROCKS"*

PolyFormer



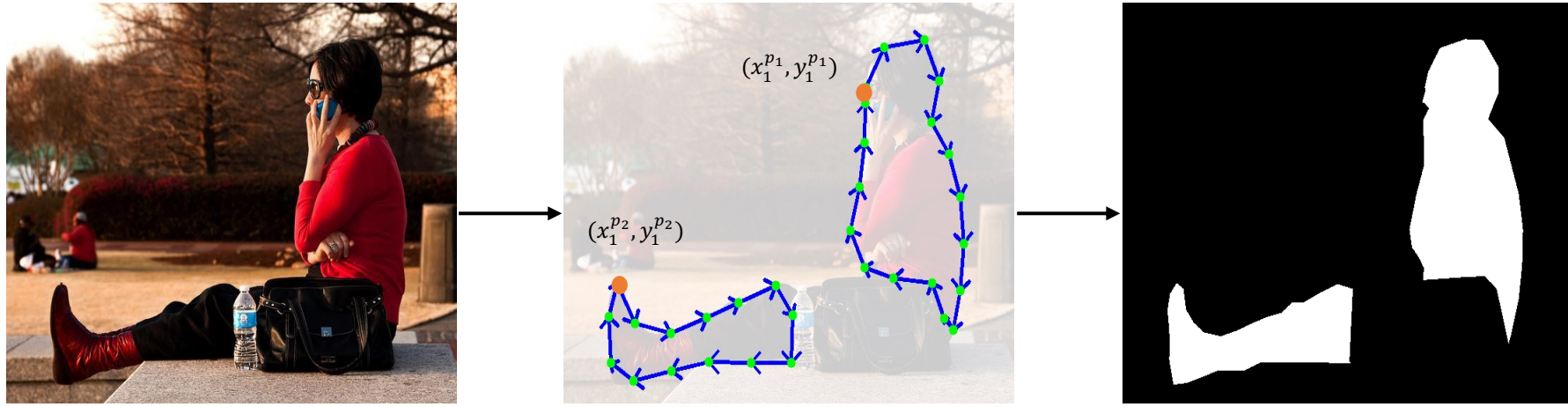
# Model Architecture





# Target Sequence Generation

- Polygon ordering
  - Start from top-left
  - Clockwise direction

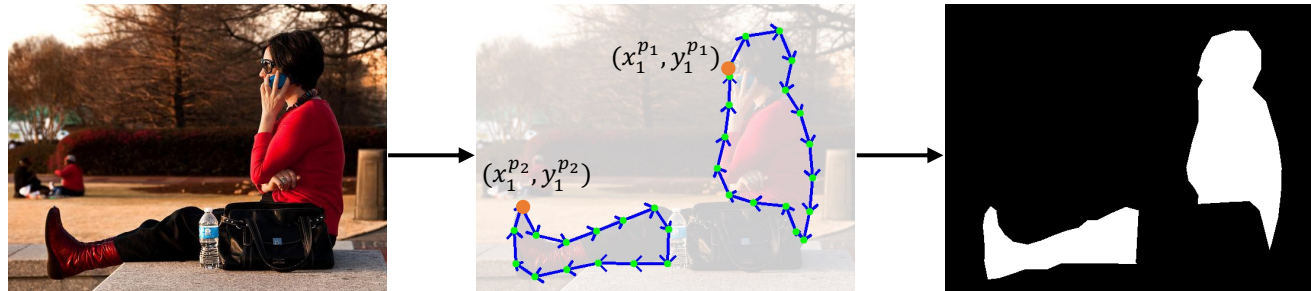
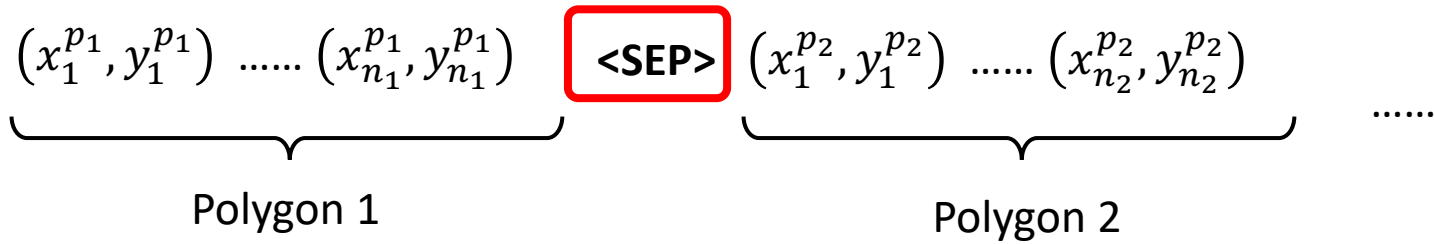


$(x_1, y_1), (x_2, y_2), \dots, (x_K, y_K)$



# Target Sequence Generation

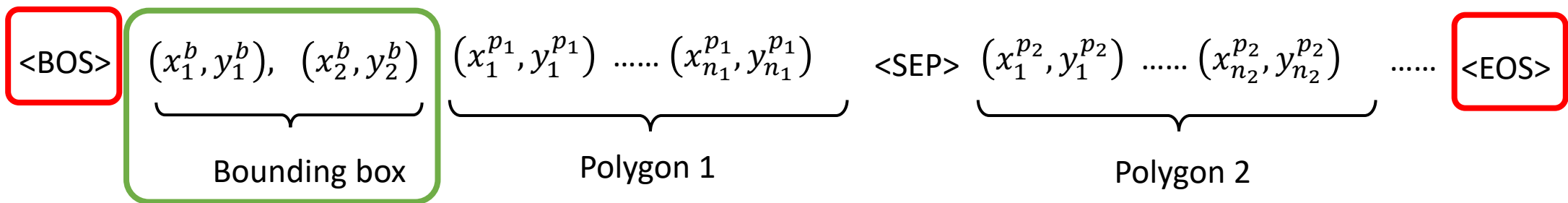
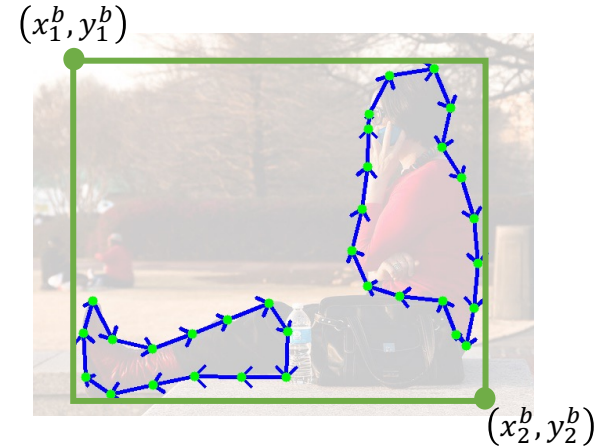
- Multi-polygon case
  - Separator token <SEP>





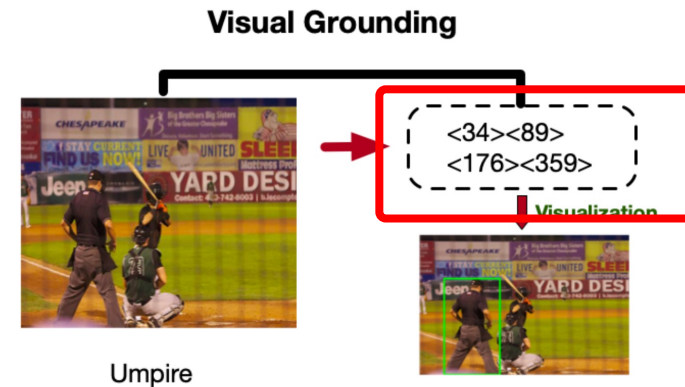
# Target Sequence Generation

- Unified sequence with bounding box
  - Bounding box:  $(x_1^b, y_1^b), (x_2^b, y_2^b)$
- Final target sequence:



# Regression-based Decoder

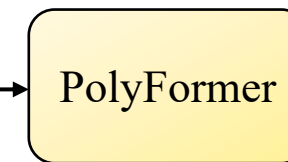
- Previous Seq2Seq framework: Coordinate prediction as a **classification task**
  - *Continuous* coordinates => *discrete* bins
  - quantization error
  - Inaccurate supervision
- PolyFormer: geometric localization as a **regression task**
  - Directly predict floating-point coordinate
  - No quantization error
  - Accurate localization



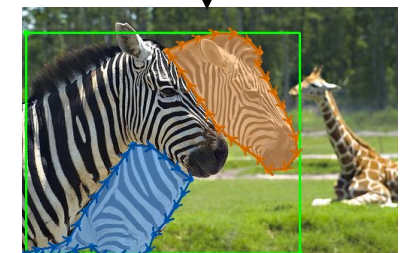
OFA [Wang et al, ICML2022]



“the second zebra from the front”



- $(x_1^b, y_1^b) = (0.5, 34.8)$
- $(x_2^b, y_2^b) = (369.0, 333.0)$
- $(x_1^{p_1}, y_1^{p_1}) = (192.7, 58.1)$
- ...
- $(x_1^{p_2}, y_1^{p_2}) = (147.5, 183.8)$
- ...

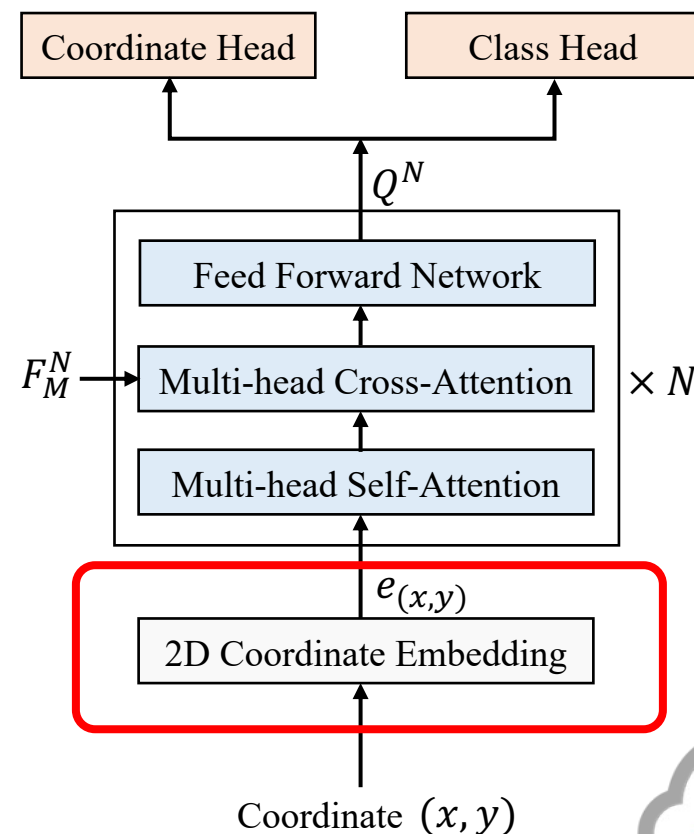
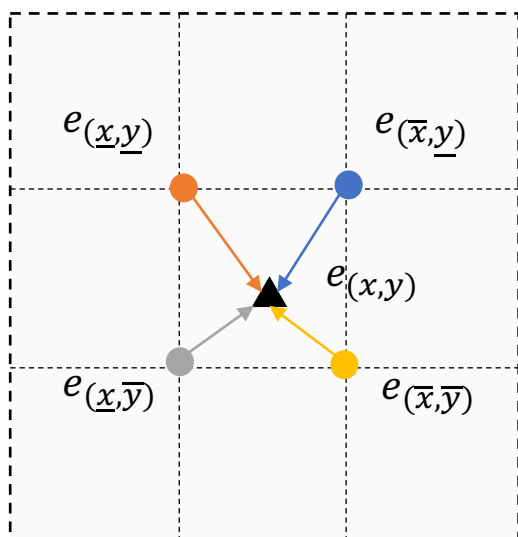


# Regression-based Transformer Decoder

- 2D Coordinate Embedding

$$e_{(x,y)} = (\bar{x} - x)(\bar{y} - y) \cdot e_{(\underline{x},\underline{y})} + (x - \underline{x})(\bar{y} - y) \cdot e_{(\bar{x},\underline{y})} + (\bar{x} - x)(y - \underline{y}) \cdot e_{(\underline{x},\bar{y})} + (x - \underline{x})(y - \underline{y}) \cdot e_{(\bar{x},\bar{y})}$$

2D Coordinate Embedding Codebook  $\mathcal{D} \in \mathbb{R}^{B_H \times B_W \times C_e}$



# Regression-based Transformer Decoder

- Prediction Heads

- **Coordinate head**

- 3-layer feed-forward network (FFN)

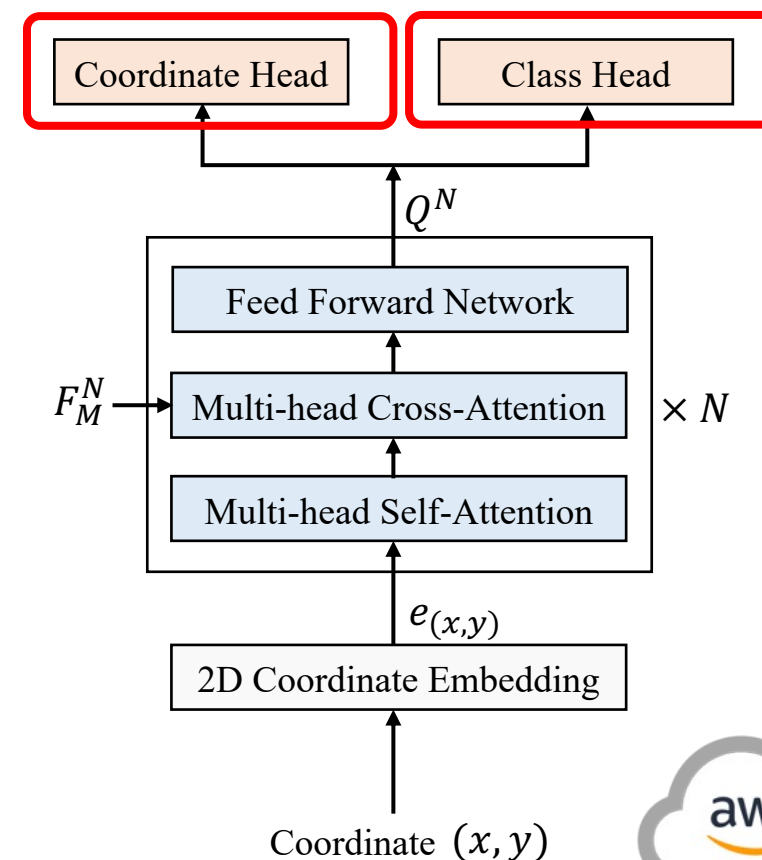
$$(\hat{x}, \hat{y}) = \text{Sigmoid}(\text{FFN}(Q^N)).$$

- **Class head**

- Linear classification layer

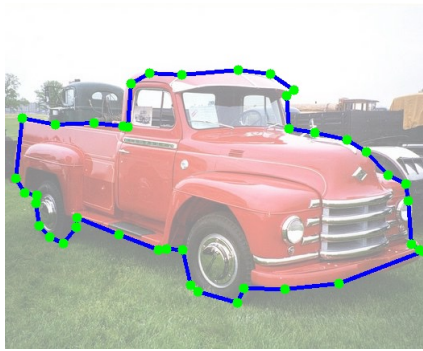
$$\hat{p} = W_c Q^N + b_c,$$

- Separator token <SEP>, coordinate token <COO>, end-of-sequence token <EOS>



# Training: Polygon Augmentation

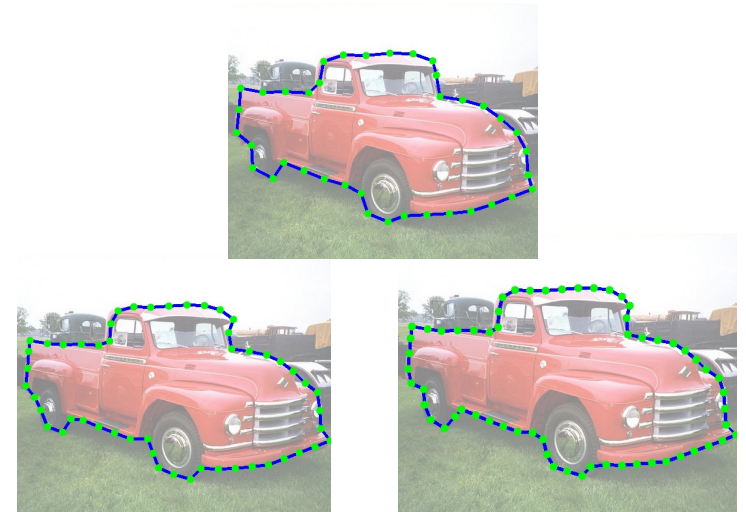
**polygons at different  
levels of granularity**



(a) Original polygon



(b) Interpolated contour



(c) Sampled polygons









# Two stage training

- Pre-train on REC task
  - Visual Genome, RefCOCO, RefCOCO+, RefCOCOg datasets, and Flickr entities
  - ~6M distinct language expressions and 164k images in the training set.
- Finetuning on REC + RIS task on RefCOCO, RefCOCO+, RefCOCOg datasets



# Referring image segmentation results

	Method	Visual Backbone	Text Encoder	RefCOCO			RefCOCO+			RefCOCOg	
				val	test A	test B	val	test A	test B	val	test
oIoU	STEP [7]	RN101	Bi-LSTM	60.04	63.46	57.97	48.19	52.33	40.41	-	-
	BRINet [29]	RN101	LSTM	60.98	62.99	59.21	48.17	52.32	42.11	-	-
	CMPC [30]	RN101	LSTM	61.36	64.53	59.64	49.56	53.44	43.23	-	-
	LSCM [31]	RN101	LSTM	61.47	64.99	59.55	49.34	53.12	43.50	-	-
	CMPC+ [49]	RN101	LSTM	62.47	65.08	60.82	50.25	54.04	43.47	-	-
	MCN [57]	DN53	Bi-GRU	62.44	64.20	59.71	50.62	54.99	44.69	49.22	49.40
	EFN [20]	WRN101	Bi-GRU	62.76	65.69	59.67	51.50	55.24	43.01	-	-
	BUSNet [81]	RN101	Self-Att	63.27	66.41	61.39	51.76	56.87	44.13	-	-
	CGAN [56]	DN53	Bi-GRU	64.86	68.04	62.07	51.03	55.51	44.06	51.01	51.69
	LTS [33]	DN53	Bi-GRU	65.43	67.76	63.08	54.21	58.32	48.02	54.40	54.25
	ReSTR [37]	ViT-B	Transformer	67.22	69.30	64.45	55.78	60.44	48.27	-	-
	PolyFormer-B	Swin-B	BERT-base	74.82	76.64	71.06	67.64	72.89	59.33	67.76	69.05
<b>PolyFormer-L</b>	Swin-L	BERT-base	<b>75.96</b>	<b>78.29</b>	<b>73.25</b>	<b>69.33</b>	<b>74.56</b>	<b>61.87</b>	<b>69.20</b>	<b>70.19</b>	
mIoU	VLT [19]	DN53	Bi-GRU	65.65	68.29	62.73	55.50	59.20	49.36	52.99	56.65
	CRIS [76]	RN101	GPT-2	70.47	73.18	66.10	62.27	68.06	53.68	59.87	60.36
	SeqTR [92]	DN53	Bi-GRU	71.70	73.31	69.82	63.04	66.73	58.97	64.69	65.74
	RefTr [42]	RN101	BERT-base	74.34	76.77	70.87	66.75	70.58	59.40	66.63	67.39
	LAVT [84]	Swin-B	BERT-base	74.46	76.89	70.94	65.81	70.97	59.23	63.34	63.62
	PolyFormer-B	Swin-B	BERT-base	75.96	77.09	73.22	70.65	74.51	64.64	69.36	69.88
	<b>PolyFormer-L</b>	Swin-L	BERT-base	<b>76.94</b>	<b>78.49</b>	<b>74.83</b>	<b>72.15</b>	<b>75.71</b>	<b>66.73</b>	<b>71.15</b>	<b>71.17</b>

PolyFormer-B  
outperforms previous  
methods on each split  
of the three datasets

Table 1. Comparison with the state-of-the-art methods on three referring image segmentation benchmarks. RN101 denotes ResNet-101 [25], WRN101 refers to Wide ResNet-101 [88], and DN53 denotes Darknet-53 [65].



# Referring image segmentation results

Method		Visual Backbone	Text Encoder	RefCOCO			RefCOCO+			RefCOCOg	
				val	test A	test B	val	test A	test B	val	test
oIoU	STEP [7]	RN101	Bi-LSTM	60.04	63.46	57.97	48.19	52.33	40.41	-	-
	BRINet [29]	RN101	LSTM	60.98	62.99	59.21	48.17	52.32	42.11	-	-
	CMPC [30]	RN101	LSTM	61.36	64.53	59.64	49.56	53.44	43.23	-	-
	LSCM [31]	RN101	LSTM	61.47	64.99	59.55	49.34	53.12	43.50	-	-
	CMPC+ [49]	RN101	LSTM	62.47	65.08	60.82	50.25	54.04	43.47	-	-
	MCN [57]	DN53	Bi-GRU	62.44	64.20	59.71	50.62	54.99	44.69	49.22	49.40
	EFN [20]	WRN101	Bi-GRU	62.76	65.69	59.67	51.50	55.24	43.01	-	-
	BUSNet [81]	RN101	Self-Att	63.27	66.41	61.39	51.76	56.87	44.13	-	-
	CGAN [56]	DN53	Bi-GRU	64.86	68.04	62.07	51.03	55.51	44.06	51.01	51.69
	LTS [33]	DN53	Bi-GRU	65.43	67.76	63.08	54.21	58.32	48.02	54.40	54.25
	ReSTR [37]	ViT-B	Transformer	67.22	69.30	64.45	55.78	60.44	48.27	-	-
	PolyFormer-B	Swin-B	BERT-base	74.82	76.64	71.06	67.64	72.89	59.33	67.76	69.05
	<b>PolyFormer-L</b>	Swin-L	BERT-base	<b>75.96</b>	<b>78.29</b>	<b>73.25</b>	<b>69.33</b>	<b>74.56</b>	<b>61.87</b>	<b>69.20</b>	<b>70.19</b>
mIoU	VLTr [19]	DN53	Bi-GRU	65.65	68.29	62.73	55.50	59.20	49.36	52.99	56.65
	CRIS [76]	RN101	GPT-2	70.47	73.18	66.10	62.27	68.06	53.68	59.87	60.36
	SeqTR [92]	DN53	Bi-GRU	71.70	73.31	69.82	63.04	66.73	58.97	64.69	65.74
	RefTr [42]	RN101	BERT-base	74.34	76.77	70.87	66.75	70.58	59.40	66.63	67.39
	LAVT [84]	Swin-B	BERT-base	74.46	76.89	70.94	65.81	70.97	59.23	63.34	63.62
		PolyFormer-B	Swin-B	BERT-base	75.96	77.09	73.22	70.65	74.51	64.64	69.36
	<b>PolyFormer-L</b>	Swin-L	BERT-base	<b>76.94</b>	<b>78.49</b>	<b>74.83</b>	<b>72.15</b>	<b>75.71</b>	<b>66.73</b>	<b>71.15</b>	<b>71.17</b>

+3.9%, 3.93%, 5.24%  
mIoU on challenging  
RefCOCO+

Table 1. Comparison with the state-of-the-art methods on three referring image segmentation benchmarks. RN101 denotes ResNet 101 [25], WRN101 refers to Wide ResNet-101 [88], and DN53 denotes Darknet-53 [65].



# Referring image segmentation results

	Method	Visual Backbone	Text Encoder	RefCOCO			RefCOCO+			RefCOCOg	
				val	test A	test B	val	test A	test B	val	test
oIoU	STEP [7]	RN101	Bi-LSTM	60.04	63.46	57.97	48.19	52.33	40.41	-	-
	BRINet [29]	RN101	LSTM	60.98	62.99	59.21	48.17	52.32	42.11	-	-
	CMPC [30]	RN101	LSTM	61.36	64.53	59.64	49.56	53.44	43.23	-	-
	LSCM [31]	RN101	LSTM	61.47	64.99	59.55	49.34	53.12	43.50	-	-
	CMPC+ [49]	RN101	LSTM	62.47	65.08	60.82	50.25	54.04	43.47	-	-
	MCN [57]	DN53	Bi-GRU	62.44	64.20	59.71	50.62	54.99	44.69	49.22	49.40
	EFN [20]	WRN101	Bi-GRU	62.76	65.69	59.67	51.50	55.24	43.01	-	-
	BUSNet [81]	RN101	Self-Att	63.27	66.41	61.39	51.76	56.87	44.13	-	-
	CGAN [56]	DN53	Bi-GRU	64.86	68.04	62.07	51.03	55.51	44.06	51.01	51.69
	LTS [33]	DN53	Bi-GRU	65.43	67.76	63.08	54.21	58.32	48.02	54.40	54.25
	ReSTR [37]	ViT-B	Transformer	67.22	69.30	64.45	55.78	60.44	48.27	-	-
	PolyFormer-B	Swin-B	BERT-base	74.82	76.64	71.06	67.64	72.89	59.33	67.76	69.05
	<b>PolyFormer-L</b>	Swin-L	BERT-base	<b>75.96</b>	<b>78.29</b>	<b>73.25</b>	<b>69.33</b>	<b>74.56</b>	<b>61.87</b>	<b>69.20</b>	<b>70.19</b>
mIoU	VLT [19]	DN53	Bi-GRU	65.65	68.29	62.73	55.50	59.20	49.36	52.99	56.65
	CRIS [76]	RN101	GPT-2	70.47	73.18	66.10	62.27	68.06	53.68	59.87	60.36
	SeqTR [92]	DN53	Bi-GRU	71.70	73.31	69.82	63.04	66.73	58.97	64.69	65.74
	RefTr [42]	RN101	BERT-base	74.34	76.77	70.87	66.75	70.58	59.40	66.63	67.39
	LAVT [84]	Swin-B	BERT-base	74.46	76.89	70.94	65.81	70.97	59.23	63.34	63.62
		PolyFormer-B	Swin-B	BERT-base	75.96	77.09	73.22	70.65	74.51	64.64	69.36
	<b>PolyFormer-L</b>	Swin-L	BERT-base	<b>76.94</b>	<b>78.49</b>	<b>74.83</b>	<b>72.15</b>	<b>75.71</b>	<b>66.73</b>	<b>71.15</b>	<b>71.17</b>

+2.73%, 2.49% mIoU on most challenging RefCOCOg

Table 1. Comparison with the state-of-the-art methods on three referring image segmentation benchmarks. RN101 denotes ResNet 101 [25], WRN101 refers to Wide ResNet-101 [88], and DN53 denotes Darknet-53 [65].





# Referring image segmentation results

	Method	Visual Backbone	Text Encoder	RefCOCO			RefCOCO+			RefCOCOg	
				val	test A	test B	val	test A	test B	val	test
oIoU	STEP [7]	RN101	Bi-LSTM	60.04	63.46	57.97	48.19	52.33	40.41	-	-
	BRINet [29]	RN101	LSTM	60.98	62.99	59.21	48.17	52.32	42.11	-	-
	CMPC [30]	RN101	LSTM	61.36	64.53	59.64	49.56	53.44	43.23	-	-
	LSCM [31]	RN101	LSTM	61.47	64.99	59.55	49.34	53.12	43.50	-	-
	CMPC+ [49]	RN101	LSTM	62.47	65.08	60.82	50.25	54.04	43.47	-	-
	MCN [57]	DN53	Bi-GRU	62.44	64.20	59.71	50.62	54.99	44.69	49.22	49.40
	EFN [20]	WRN101	Bi-GRU	62.76	65.69	59.67	51.50	55.24	43.01	-	-
	BUSNet [81]	RN101	Self-Att	63.27	66.41	61.39	51.76	56.87	44.13	-	-
	CGAN [56]	DN53	Bi-GRU	64.86	68.04	62.07	51.03	55.51	44.06	51.01	51.69
	LTS [33]	DN53	Bi-GRU	65.43	67.76	63.08	54.21	58.32	48.02	54.40	54.25
	ReSTR [37]	ViT-B	Transformer	67.22	69.30	64.45	55.78	60.44	48.27	-	-
	PolyFormer-B	Swin-B	BERT-base	74.82	76.64	71.06	67.64	72.89	59.33	67.76	69.05
	<b>PolyFormer-L</b>	Swin-L	BERT-base	<b>75.96</b>	<b>78.29</b>	<b>73.25</b>	<b>69.33</b>	<b>74.56</b>	<b>61.87</b>	<b>69.20</b>	<b>70.19</b>
mIoU	VLT [19]	DN53	Bi-GRU	65.65	68.29	62.73	55.50	59.20	49.36	52.99	56.65
	CRIS [76]	RN101	GPT-2	70.47	73.18	66.10	62.27	68.06	53.68	59.87	60.36
	SeqTR [92]	DN53	Bi-GRU	71.70	73.31	69.82	63.04	66.73	58.97	64.69	65.74
	RefTr [42]	RN101	BERT-base	74.34	76.77	70.87	66.75	70.58	59.40	66.63	67.39
	LAVT [84]	Swin-B	BERT-base	74.46	76.89	70.94	65.81	70.97	59.23	63.34	63.62
		PolyFormer-B	Swin-B	BERT-base	75.96	77.09	73.22	70.65	74.51	64.64	69.36
	<b>PolyFormer-L</b>	Swin-L	BERT-base	<b>76.94</b>	<b>78.49</b>	<b>74.83</b>	<b>72.15</b>	<b>75.71</b>	<b>66.73</b>	<b>71.15</b>	<b>71.17</b>

PolyFormer-L vs. B:  
+1~2 points

Table 1. Comparison with the state-of-the-art methods on three referring image segmentation benchmarks. RN101 denotes ResNet-101 [25], WRN101 refers to Wide ResNet-101 [88], and DN53 denotes Darknet-53 [65].



# Zero-shot Transfer to Referring Video Object Segmentation

Method	Visual Backbone	$\mathcal{J}\&\mathcal{F}$	$\mathcal{J}$	$\mathcal{F}$
CMSA+RNN [85]	ResNet-50	40.2	36.9	43.5
URVOS [70]	ResNet-50	51.5	47.3	56.0
CITD [44]	ResNet-101	56.4	54.8	58.1
ReferFormer [78]	Swin-L	60.5	57.6	63.4
<b>ReferFormer [78]</b>	<b>Video-Swin-B</b>	<b>61.1</b>	<b>58.1</b>	<b>64.1</b>
PolyFormer-B <sup>†</sup>	Swin-B	60.9	56.6	65.2
<b>PolyFormer-L<sup>†</sup></b>	<b>Swin-L</b>	<b>61.5</b>	<b>57.2</b>	<b>65.8</b>

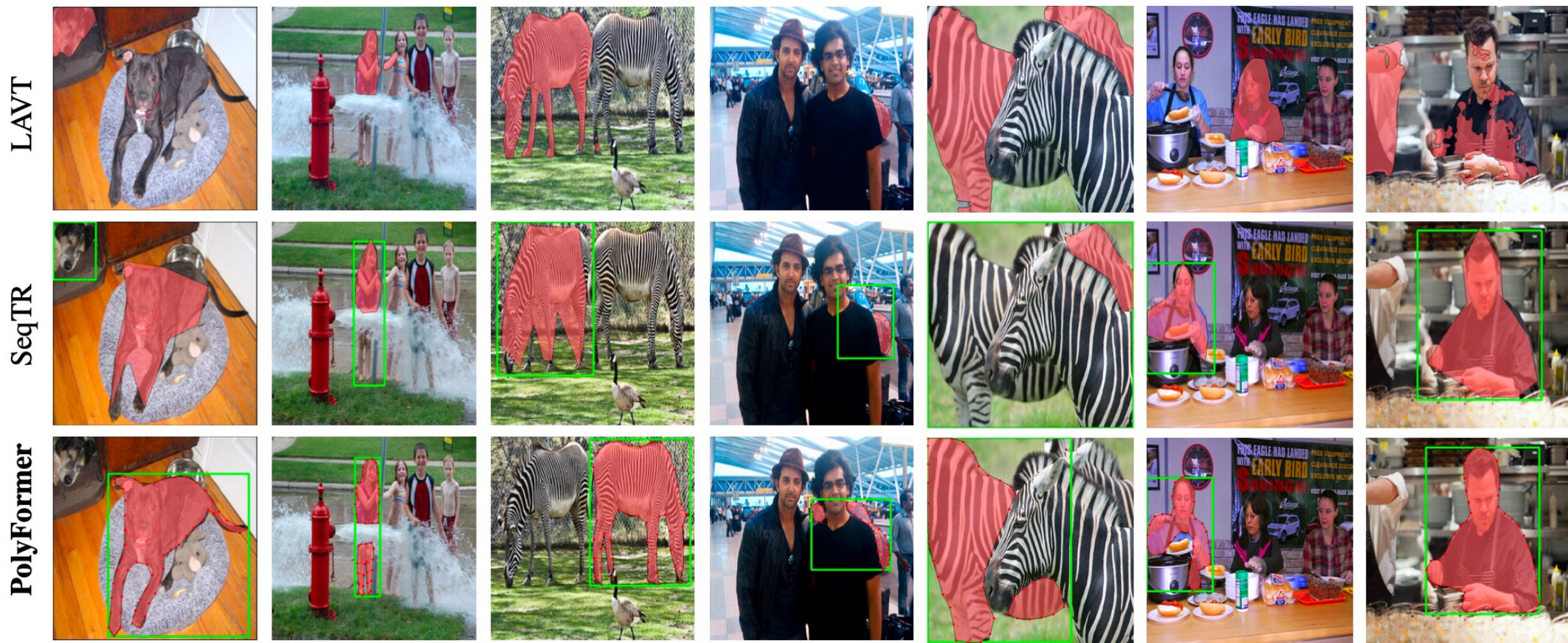
Best J&F w/o training on video

Table 3. Comparison with the state-of-the-art methods on Ref-DAVIS17. † means our model is trained on image datasets only. ReferFormer is trained on both image and video datasets.





# Visualization Results on RefCOCOg



(a) “a dark grey dog on a light grey round bed wearing a red collar”

(b) “girl in purple”

(c) “zebra eating grass with a goose in front of it”

(d) “a black car parked at a transportation terminal”

(e) “a zebra with its head not visible but much of its body able to be seen”

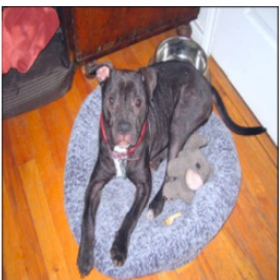

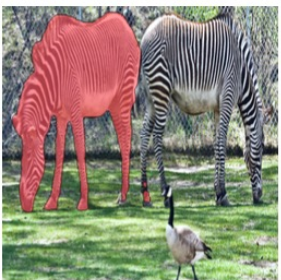









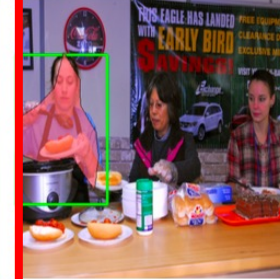

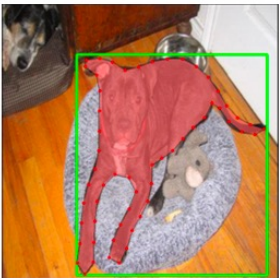

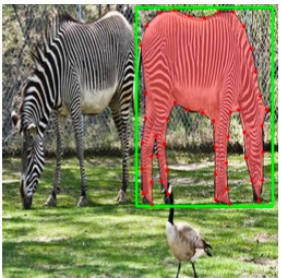




(f) “a girl was cooking the food and serving”

(g) “a man wearing a black shirt and a black and white striped apron stirring something in a metal container”



# Visualization Results on RefCOCOg

complex language understanding

LAVT							
SeqTR							
PolyFormer							
	(a) "a dark grey dog on a light grey round bed wearing a red collar"	(b) "girl in purple"	(c) "zebra eating grass with a goose in front of it"	(d) "a black car parked at a transportation terminal"	(e) "a zebra with its head not visible but much of its body able to be seen"	(f) "a girl was cooking the food and serving"	(g) "a man wearing a black shirt and a black and white striped apron stirring something in a metal container"



# Visualization Results on RefCOCOg

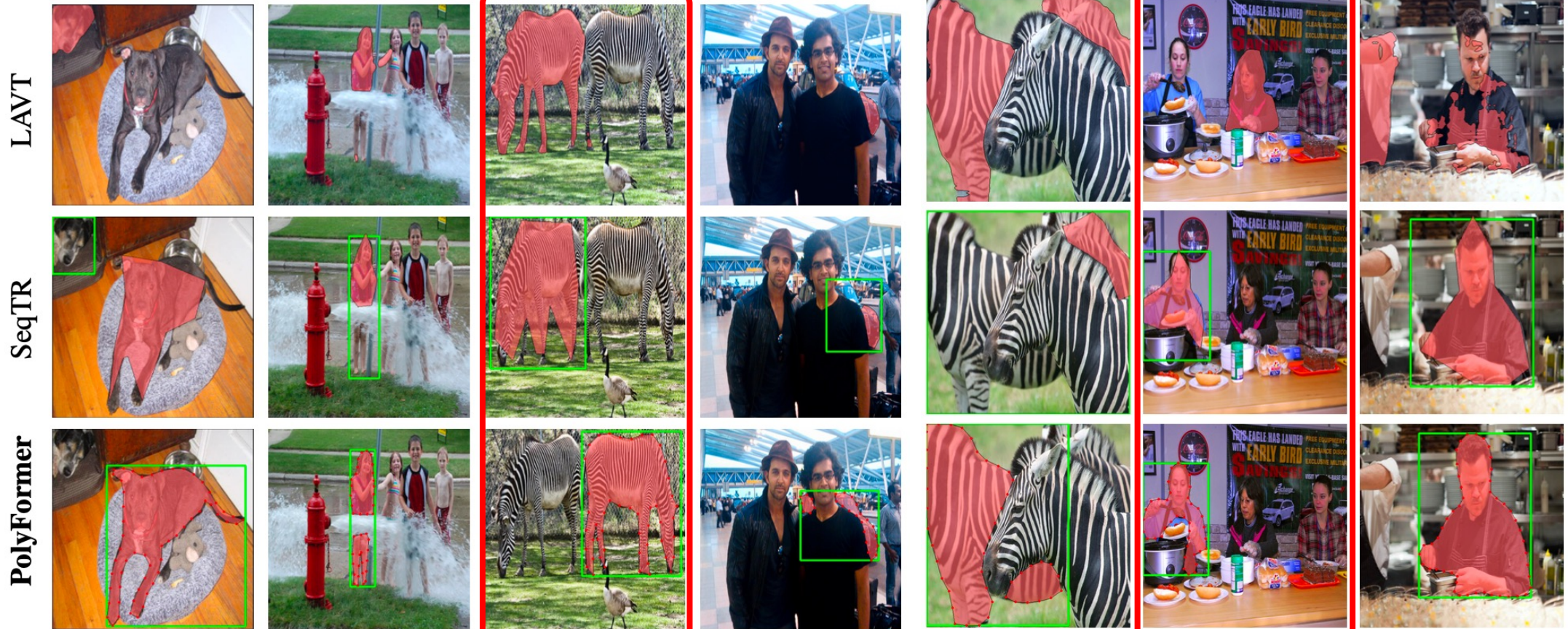
Instances with occlusion

LAVT							
SeqTR							
PolyFormer							
	(a) "a dark grey dog on a light grey round bed wearing a red collar"	(b) "girl in purple"	(c) "zebra eating grass with a goose in front of it"	(d) "a black car parked at a transportation terminal"	(e) "a zebra with its head not visible but much of its body able to be seen"	(f) "a girl was cooking the food and serving"	(g) "a man wearing a black shirt and a black and white striped apron stirring something in a metal container"



# Visualization Results on RefCOCOg

complex vision-language semantics



(a) “a dark grey dog on a light grey round bed wearing a red collar”

(b) “girl in purple”

(c) “zebra eating grass with a goose in front of it”

(d) “a black car parked at a transportation terminal”

(e) “a zebra with its head not visible but much of its body able to be seen”

(f) “a girl was cooking the food and serving”

(g) “a man wearing a black shirt and a black and white striped apron stirring something in a metal container”



# Zero-shot Evaluation on Stable Diffusion Images



(a) "A cat chef cooking fish in a fancy restaurant"  
 (b) "A chair that looks like octopus"  
 (c) "A small cabin on top of a snowy mountain in the style of Disney artstation"  
 (d) "A shiba inu puppy painted by Monet"  
 (e) "A gentleman otter in a 19th century portrait"  
 (f) "A pikachu fine-dining with a view to the Eiffel Tower"  
 (g) "A pig robot preparing a delicious meal"



# PolyFormer: Referring Image Segmentation as Sequential Polygon Generation



Jiang Liu<sup>1\*†</sup>



Hui Ding<sup>2\*</sup>



Zhaowei Cai<sup>2</sup>



Yuting Zhang<sup>2</sup>



Ravi Kumar  
Satzoda<sup>2</sup>



Vijay Mahadevan<sup>2</sup>



R. Manmatha<sup>2</sup>

<sup>1</sup> Johns Hopkins University, <sup>2</sup> AWS AI Labs, \*Equal Contribution, †Work done during internship at AWS AI Labs

<https://polyformer.github.io/>

