

# CoupAlign: Coupling Word-Pixel with Sentence-Mask Alignments for Referring Image Segmentation

---





Zicheng Zhang<sup>1\*</sup> Yi Zhu<sup>2\*</sup> Jianzhuang Liu<sup>2</sup> Xiaodan Liang<sup>3</sup> Wei Ke<sup>1†</sup>

<sup>1</sup>Xi'an Jiaotong University   <sup>2</sup>Noah's Ark Lab, Huawei   <sup>3</sup>Sun Yat-sen University

# Background

- Referring image segmentation (RIS) aims at localizing all pixels of the visual objects described by a natural language sentence.
- Main challenges:
  - how to align the given language expression with visual pixels for highlighting the target.
  - how to distinguish the target from similar objects.

**Input:** image + texts

-  white double decker bus
-  gray small car
-  the old man with T-shirt
-  woman with pink dress

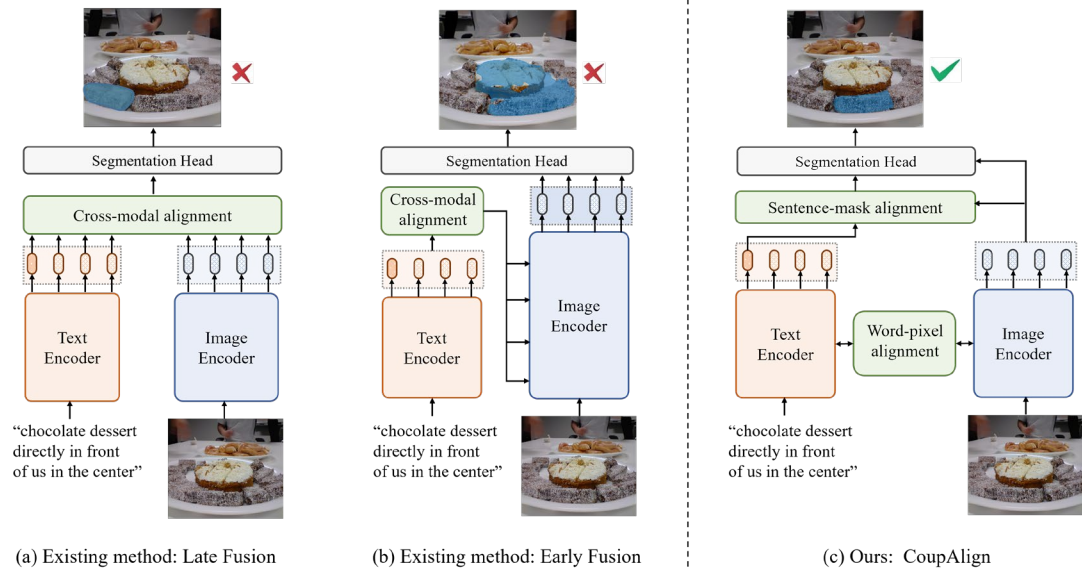


**Output:** instance masks



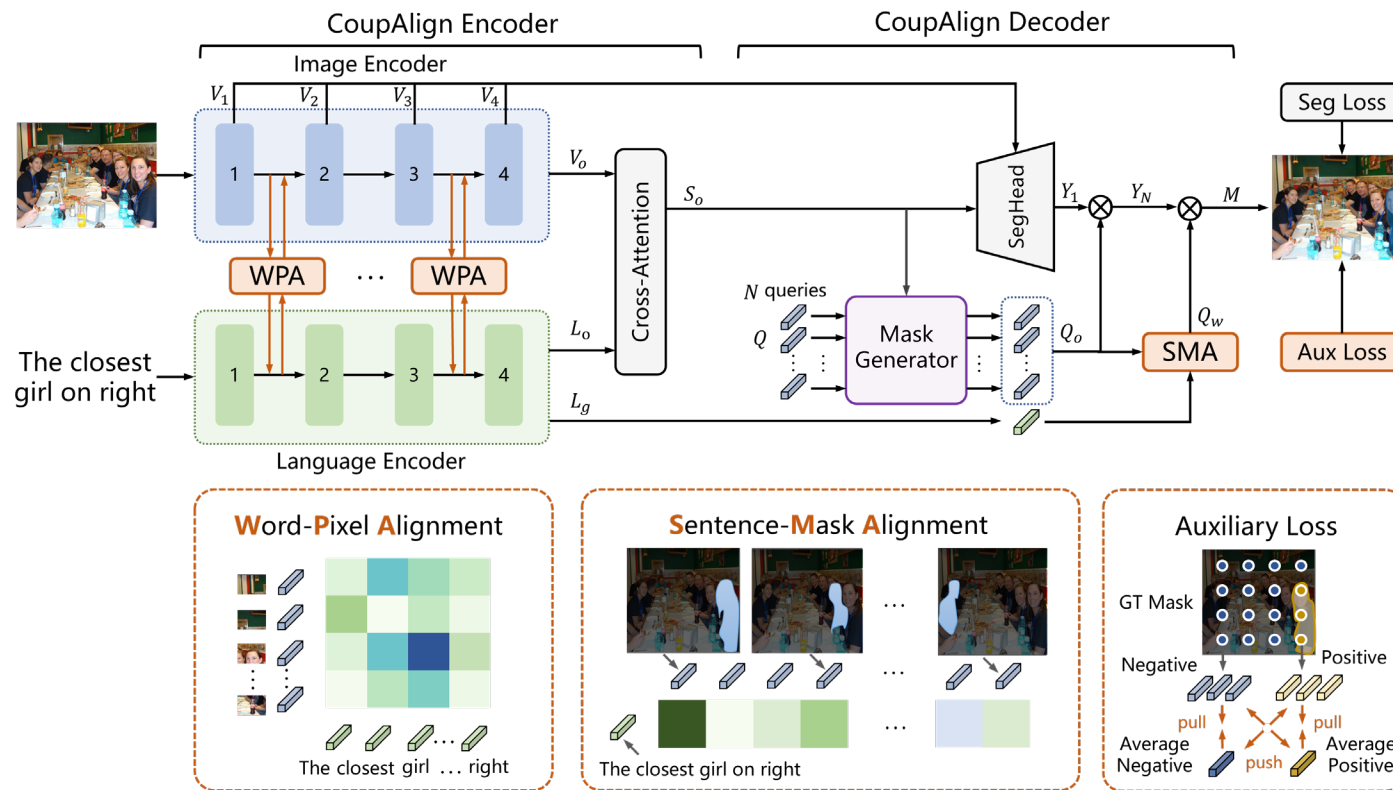
# Motivation

- Previous methods
  - only consider late fusion for vision and language features.
  - or consider early fusion but only have word-pixel level alignment.
- Our methods
  - adopts word-pixel level alignment in the early fusion stage.
  - and use sentence-mask level alignment to enhance the fused features in the late fusion stage.



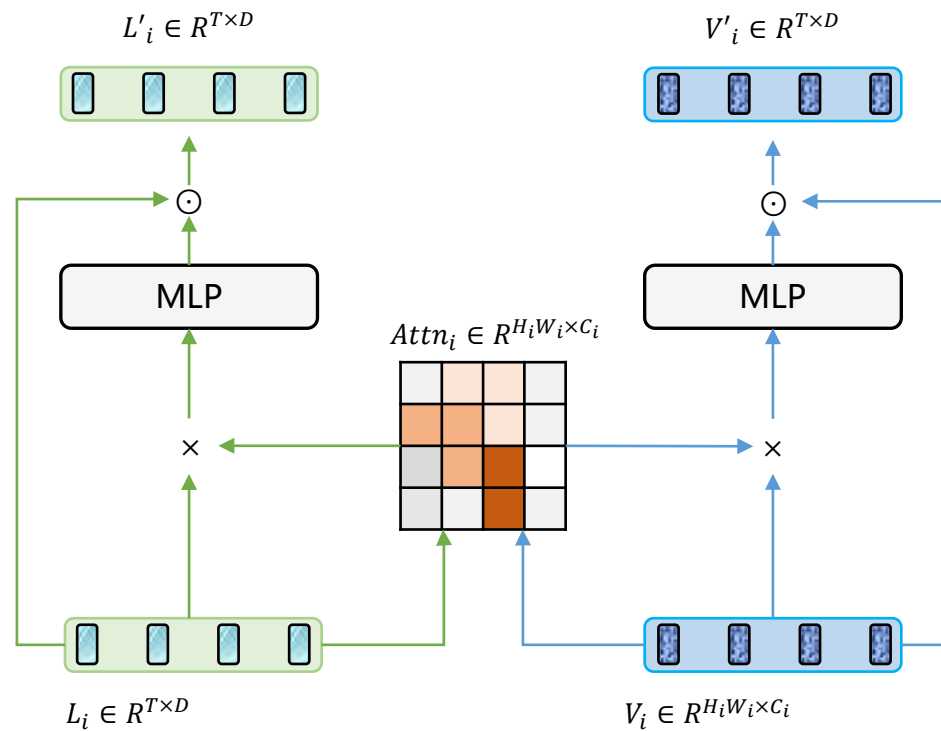
## ➤ CoupAlign Architecture

- WPA module enables cross-model interactions at each encoder stage.
- Based on the aligned feature  $S_o$ , the mask generator produces N mask embeddings.
- SMA module weights  $Q_o$  using  $L_g$  and projects the mask signals back to  $Y_N$ .

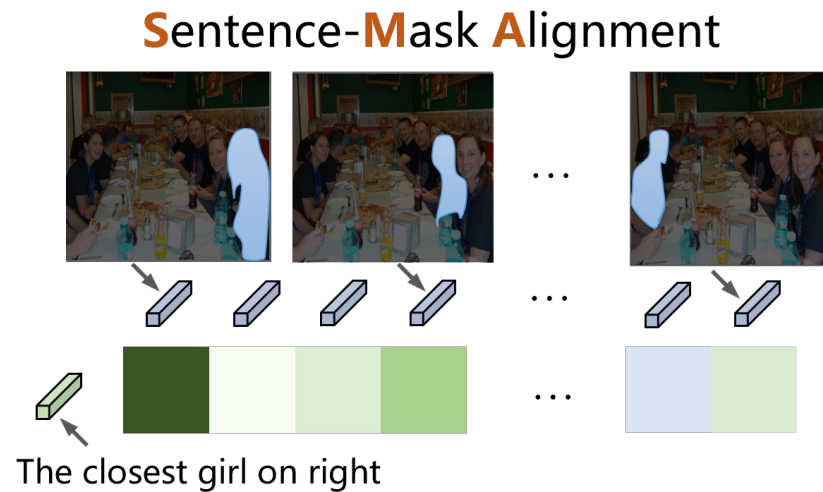


# Method

- Word-Pixel Alignment
  - we use cross attention to align word tokens and pixel tokens
  - and design a gate to control the fused information flow.

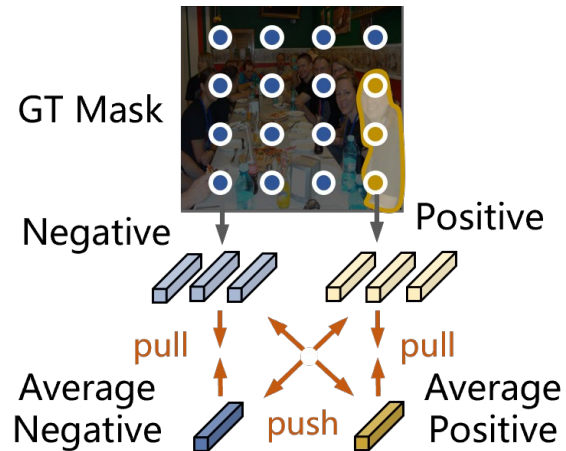


- Sentence-Mask Alignment
  - we use global language feature and mask embeddings to compute attention weights.
  - Then we use the weights to aggregate the proposals to the final mask prediction.



➤ Auxiliary loss

- We adopt contrastive loss to enhance the ability of distinguishing the object from the background.



$$L_{P2N}^{\text{Aux}} = -\frac{1}{|\mathcal{P}|} \sum_{y_i^+ \in \mathcal{P}} \frac{\exp(y_i^+ \cdot \hat{y}^+ / \tau)}{\exp(y_i^+ \cdot \hat{y}^+ / \tau) + \sum_{y_k^- \in \mathcal{N}} \exp(y_i^+ \cdot y_k^- / \tau)},$$

$$L_{N2P}^{\text{Aux}} = -\frac{1}{|\mathcal{N}|} \sum_{y_i^- \in \mathcal{N}} \frac{\exp(y_i^- \cdot \hat{y}^- / \tau)}{\exp(y_i^- \cdot \hat{y}^- / \tau) + \sum_{y_k^+ \in \mathcal{P}} \exp(y_i^- \cdot y_k^+ / \tau)},$$

$$L_j^{\text{Aux}} = L_{P2N}^{\text{Aux}} + L_{N2P}^{\text{Aux}},$$

# Result

## ➤ Comparing with SOTA

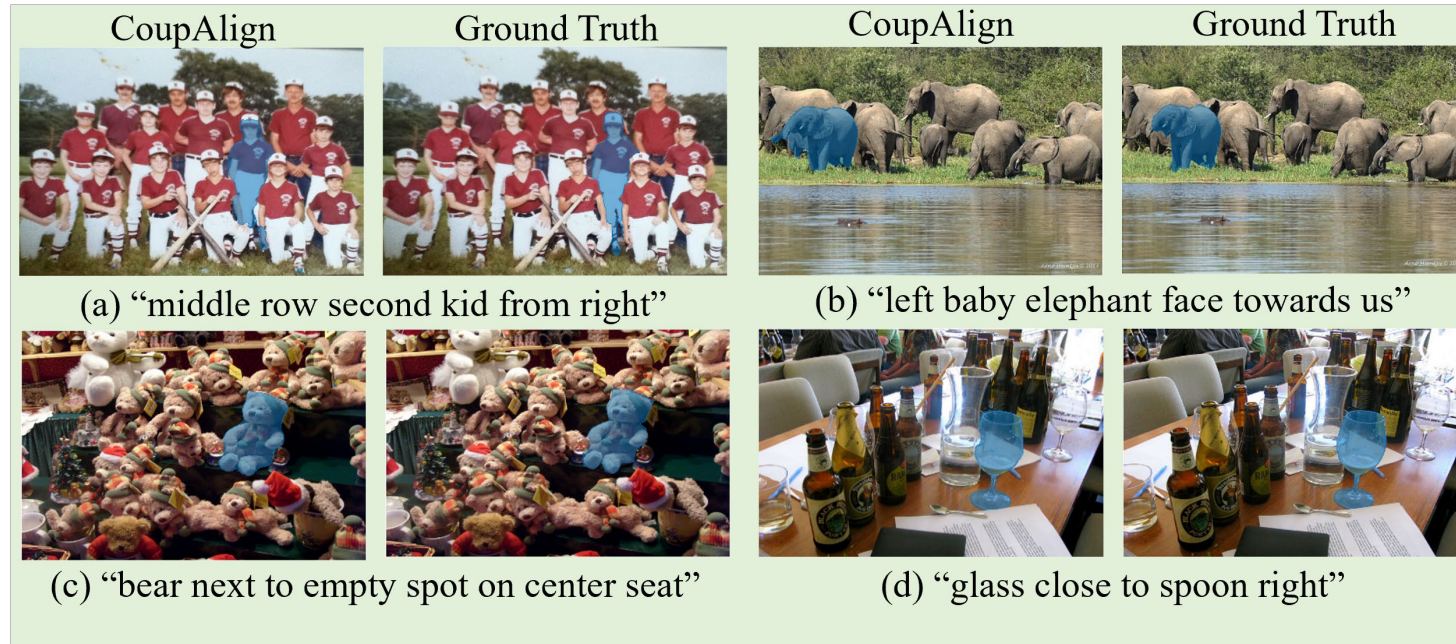
- We evaluate our method on four referring image segmentation dataset and outperform previous method.

	Backbone	RefCOCO			RefCOCO+			G-Ref		ReferIt
		val	test A	test B	val	testA	testB	val (U)	test (U)	test
MCN [31]	Darknet-53	62.44	64.20	59.71	50.62	54.99	44.69	49.22	49.40	-
BRINet [12]	ResNet-101	60.98	62.99	59.21	48.17	52.32	42.11	-	-	63.46
CMPC [13]	ResNet-101	61.36	64.53	59.64	49.56	53.44	43.23	-	-	65.53
LSCM [14]	ResNet-101	61.47	64.99	59.55	49.34	53.12	43.50	-	-	66.57
CGAN [30]	ResNet-101	64.86	68.04	62.07	51.03	55.51	44.06	51.01	51.69	-
BUSNet [43]	ResNet-101	63.27	66.41	61.39	51.76	56.87	44.13	-	-	-
EFN [11]	ResNet-101	62.76	65.69	59.67	51.50	55.24	43.01	-	-	66.70
LTS [16]	DarkNet-53	65.43	67.76	63.08	54.21	58.32	48.02	54.40	54.25	-
VLT [7]	DarkNet-53	65.65	68.29	62.73	55.50	59.20	49.36	52.99	56.65	-
ReSTR [19]	ViT-B-16	67.22	69.30	64.45	55.78	60.44	48.27	54.48	-	70.18
CRIS [41]	ResNet-101	70.47	73.18	66.10	62.27	68.08	53.68	59.87	60.36	-
LAVT [44]	Swin-B	72.73	75.82	68.79	62.14	<b>68.38</b>	55.10	61.24	62.09	-
CoupAlign (ours)	Swin-B	<b>74.70</b>	<b>77.76</b>	<b>70.58</b>	<b>62.92</b>	68.34	<b>56.69</b>	<b>62.84</b>	<b>62.22</b>	<b>73.28</b>



# Visualization

- Visualization
  - CoupAlign works well in the scenes where crowded objects have similar color and context.



# Conclusion

---

## ➤ Conclusion

- CoupAlign captures both visual and semantic coherence of pixels within the referred object, and significantly outperforms state-of-the-art RIS methods.
- Especially, CoupAlign has great ability in localizing the target from similar objects, showing great potential in segmenting natural language referred objects in real-world scenarios.

# Thank you!

---