



**VIP Lab**  
Visual Information Processing



# Finding NeMo 🐼 : Negative-mined Mosaic Augmentation for Referring Image Segmentation

Seongsu Ha<sup>1,2\*</sup>, Chaeyun Kim<sup>1\*</sup>, Donghwa Kim<sup>1\*</sup>, Junho Lee<sup>1</sup>, Sangho Lee<sup>3</sup>, Joonseok Lee<sup>1,4</sup>

Seoul National University <sup>1</sup> Twelve Labs <sup>2</sup> Ai2 <sup>3</sup> Google Research <sup>4</sup>

# Referring Image Segmentation

- Given an image and a text, RIS predicts a segmentation mask of the object referred.
- The key to RIS is to discern the referent among visually similar objects via textual cues.

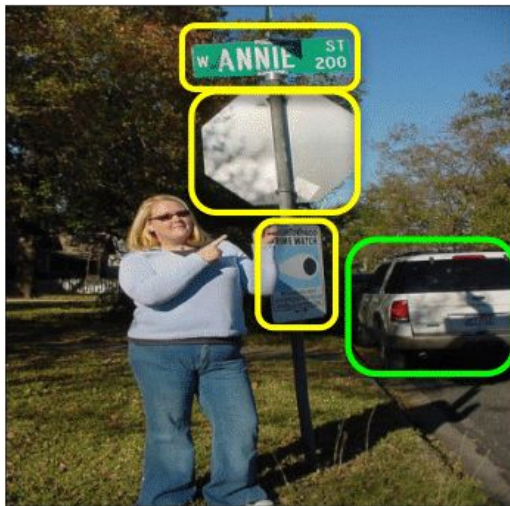
a young woman in blue shirt and striped pants sitting in the snow



a skier in an orange jacket bending over

# What makes RIS difficult?

- The difficulty of each RIS scenario can be affected by the degree of visual ambiguity in the scene given the linguistic complexity of the referring expression.



- (1) “a **sign** lettered ‘ANNIE’ between woman and SUV”
- (2) “a parked white FORD **SUV**”

# Motivation - the gap between easy & hard scenarios

- We manually pick 100 easy and hard samples depending on the number of negative objects.
- A huge performance gap exists between easy & hard examples in current models.

**Table 2:** mIoU & oIoU on 100 easy and hard samples from G-Ref UMD test set

Models	mIoU		oIoU	
	Easy	Hard	Easy	Hard
LAVT [50]	78.26	54.61	79.16	47.40
CRIS [45]	76.89	52.97	78.81	43.20
CGFormer [43]	79.86	61.22	79.95	53.27

**Fig. I:** Easy samples from G-Ref test split



*"A little girl in a blue dress"*

*"A yellow train with black trim"*

**Fig. II:** Hard samples from G-Ref test split



*"An uncooked pizza with four hotdogs"*

*"A white toothbrush with green, blue and white bristles"*

# Motivation - training data challenging enough?

- Variant Inter and even Intra-dataset grounding difficulty levels exist in training data as well.
- We ask if these samples are challenging enough to discern subtle visual and textual nuance for RIS.

RefCOCO(+)



*“right woman”  
“the woman”*

G-Ref (hard)



*“a woman getting  
her hair brushed”*

G-Ref (easy)



*“woman holding a glass  
and wearing white t-shirt”*



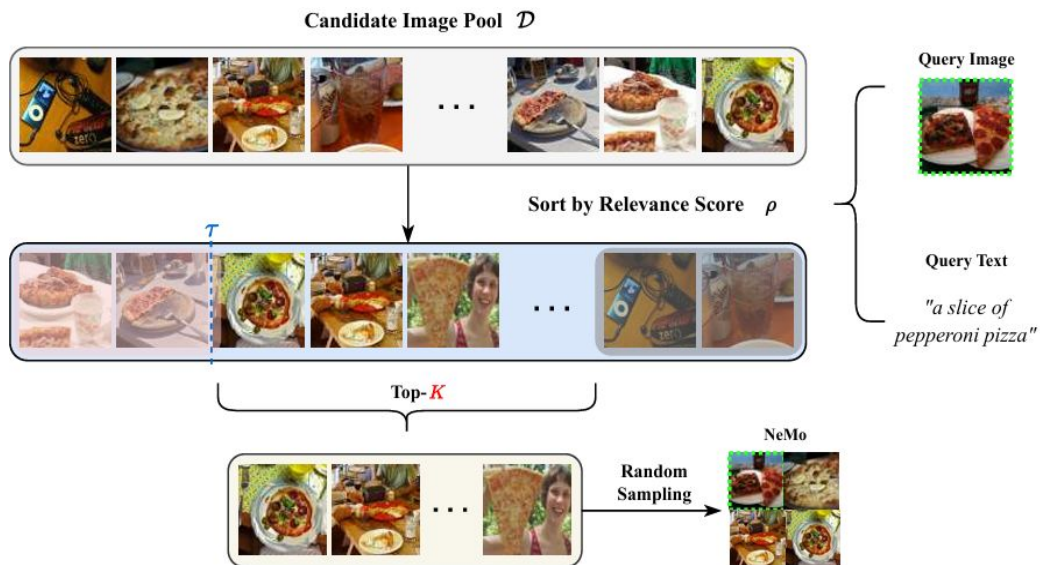
# Overall Pipeline

- **NeMo: Negative Mining + Mosaic Augmentation**
- **Filtering is necessary for the right level of ambiguity, and to avoid invalid mosaics.**



“the rightmost pizza on a paper plate”

“a man jumping with a skateboard”



# Overall Comparison

- **Overall RIS performance (oIoU) comparison w/ and w/o NeMo**

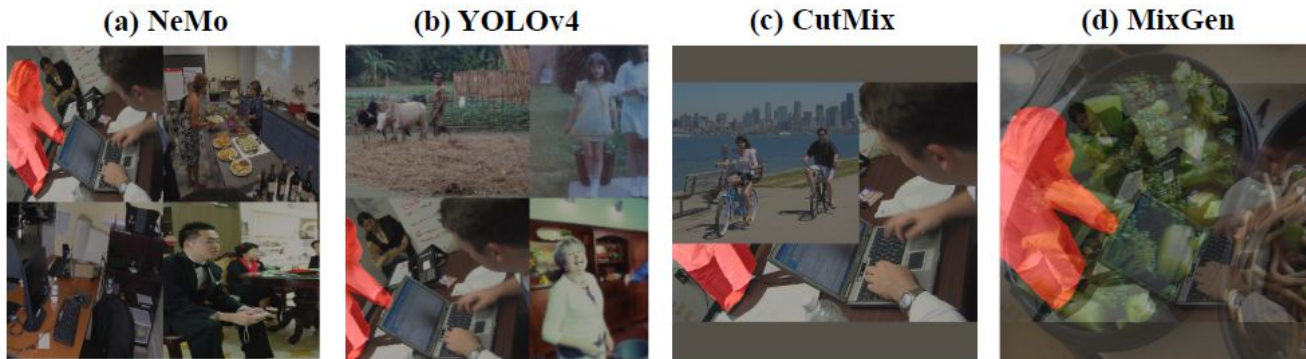
- We observe a larger performance boost on more complex datasets.
- Harder datasets benefit more because of its intricate referring expressions and visually dense scenes.

RIS model	NeMo	RefCOCO (UNC)			RefCOCO+ (UNC)			G-Ref (UMD)		GRES Val	Average Gain
		Val	TestA	TestB	Val	TestA	TestB	Val	Test		
LAVT [50]	✗	72.73	75.82	68.79	62.14	68.38	55.10	61.24	62.09	57.64	$+1.92_{\pm 2.34}$
	✓	<b>73.25</b>	<b>76.12</b>	<b>69.67</b>	<b>62.52</b>	<b>69.95</b>	<b>56.02</b>	<b>63.40</b>	<b>64.95</b>	<b>65.35</b>	
CRIS [45]	✗	66.68	70.62	59.93	56.94	64.20	46.97	55.91	58.50	54.55	$+1.73_{\pm 0.82}$
	✓	<b>68.66</b>	<b>72.82</b>	<b>63.06</b>	<b>57.94</b>	<b>65.25</b>	<b>48.41</b>	<b>58.47</b>	<b>59.07</b>	<b>56.23</b>	
ReLA [28]	✗	73.67	76.18	70.39	63.82	68.70	55.78	65.22	65.29	63.10	$+0.97_{\pm 0.82}$
	✓	<b>74.24</b>	<b>77.11</b>	<b>70.39</b>	<b>65.35</b>	<b>70.55</b>	<b>56.68</b>	<b>65.32</b>	<b>65.73</b>	<b>65.54</b>	
CGFormer [43]	✗	72.53	75.12	70.09	63.55	68.58	56.05	62.92	64.63	64.77	$+1.04_{\pm 0.67}$
	✓	<b>73.52</b>	<b>76.07</b>	<b>70.92</b>	<b>64.30</b>	<b>69.58</b>	<b>57.85</b>	<b>65.31</b>	<b>65.07</b>	<b>65.00</b>	
VPD [58]	✗	73.46	75.31	70.23	61.41	67.98	54.99	63.12	63.59	62.38	$+1.47_{\pm 0.85}$
	✓	<b>74.48</b>	<b>76.32</b>	<b>71.51</b>	<b>62.86</b>	<b>69.92</b>	<b>55.56</b>	<b>64.40</b>	<b>64.80</b>	<b>65.89</b>	
Average Gain		$+1.11_{\pm 0.79}$			$+1.21_{\pm 0.48}$			$+1.55_{\pm 0.99}$		$+3.11_{\pm 2.83}$	

**Table 3:** Overall RIS performance (in oIoU) comparison with and without NeMo



# Comparison to other augmentation methods



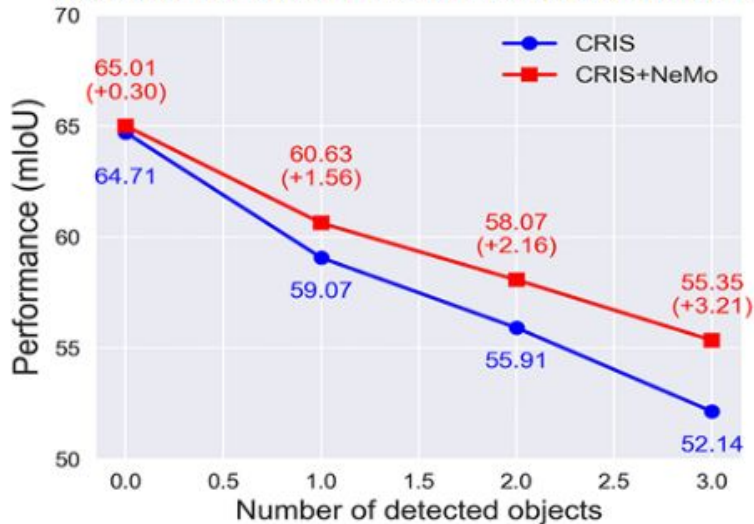
Query : A woman in a white shirt looking down at a laptop

Augmentation Method	oIoU		Prec (Val)	
	Val	Test	0.5	0.7
CRIS	55.91	58.50	67.95	54.84
+YOLOv4 [3]	56.22	58.55	66.94	53.54
+CutMix [56]	56.50	58.34	66.63	53.11
+MixGen [13]	53.62	55.85	64.37	51.28
+NeMo (Ours)	<b>58.47</b>	<b>59.07</b>	<b>70.01</b>	<b>56.60</b>

# Detailed Analysis (1)

- **Performance on Visually Challenging Scenarios**
  - better in challenging cases with more negative objects.
- **Performance w.r.t Query Complexity**
  - robust at sentence lengths, even with longer complex ones.

(1) Performance on Visually Challenging Scenarios



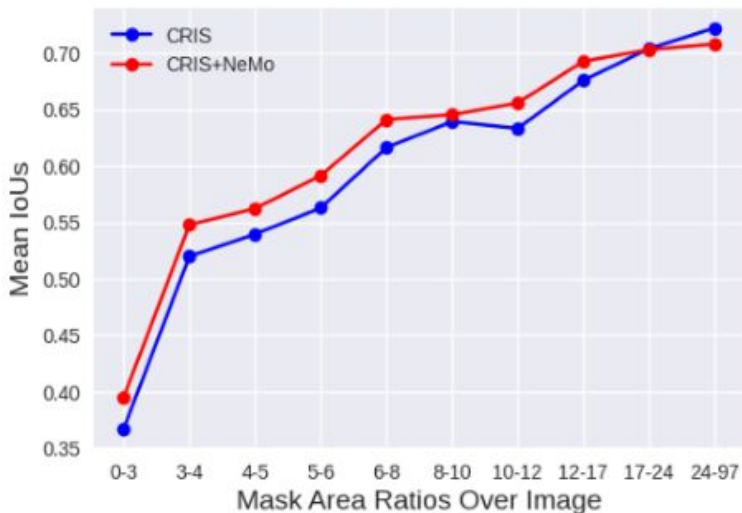
(2) Performance w.r.t Query Complexity

RIS model	NeMo	Length of $T$			
		1-5	6-7	8-10	11-20
LAVT [50]	X ✓	63.95 <b>66.50</b>	63.46 <b>65.39</b>	63.03 <b>64.40</b>	63.00 <b>64.72</b>
CRIS [45]	X ✓	58.91 <b>60.77</b>	56.41 <b>57.17</b>	55.29 <b>57.05</b>	57.33 <b>58.35</b>
ReLA [28]	X ✓	<b>66.67</b> 66.63	64.95 <b>65.00</b>	<b>63.82</b> 63.75	65.95 <b>67.26</b>
CGFormer [43]	X ✓	65.85 <b>66.30</b>	65.12 <b>65.44</b>	<b>64.33</b> 63.98	63.87 <b>64.98</b>
VPD [58]	X ✓	<b>67.53</b> 66.30	66.12 <b>66.86</b>	65.49 <b>67.33</b>	67.44 <b>68.12</b>

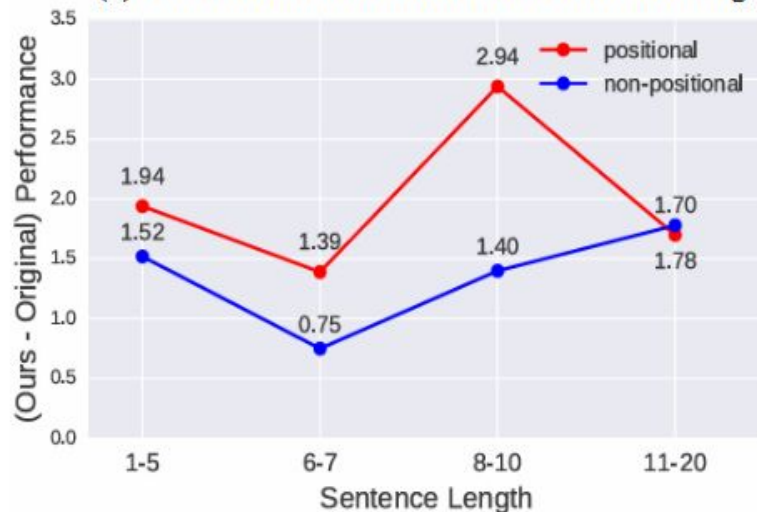
# Detailed Analysis (2)

- **Robustness on Object Scale**
  - better in most object sizes, especially for smaller objects.
- **Enhancement on Positional Understanding**
  - better at positional keywords, even in long and complex queries.

(3) Robustness on Object Scale (G-Ref val set)



(4) Enhancement on Positional Understanding



# Qualitative Analysis (1)

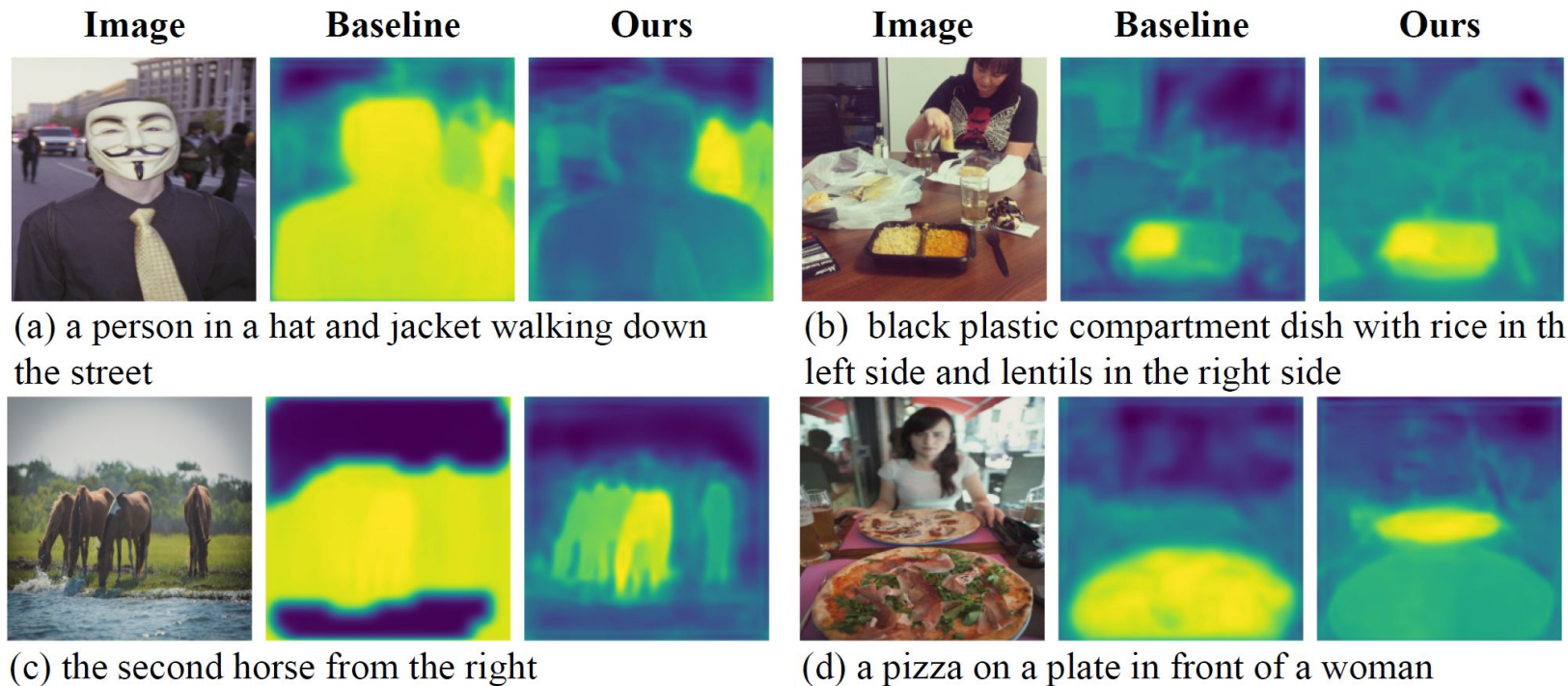


Fig. 9: Visualization of activation maps with and without NeMo on CRIS

# Qualitative Analysis (2)

*“a lady in a black dress cuts a wedding cake with her new husband”*



*“a slice of cheese cake at the top of the fork”*



(a) *“a man with a white cap and brown shirt standing next to an elephant”*



(b) *“giraffe holding head highest”*



(c) *“the front edge of a tan scooter with a carrying container on it”*



**Image**

**Baseline**

**Ours**

**Ground Truth**

Figure 7. Visualization of results after augmentation on CRIS [45].

# Summary

---

- We introduce NeMo, Negative-mined Mosaic Augmentation, a simple but powerful labor-free data augmentation method for Referring Image Segmentation.
- NeMo involves a systematic way to tune the dataset difficulty by generating training examples at a properly controlled difficulty.
- NeMo brings consistent IoU improvement over various state-of-the-art RIS models on multiple datasets.
- NeMo enhances both visual and textual understanding capabilities for segmenting the right target.