



# 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>, and Joonseok Lee<sup>1,4\*\*</sup>

<sup>1</sup> Seoul National University, Seoul, Korea

<sup>2</sup> Twelve Labs, Seoul, Korea

<sup>3</sup> Allen Institute for AI, Seattle, Washington, USA

<sup>4</sup> Google Research, Mountain View, California, USA

`mars@twelvelabs.io`, `{golddohyun, ehd9712, joon2003}@snu.ac.kr`,  
`sanghol@allenai.org`, `joonseok@snu.ac.kr`

**Abstract.** Referring Image Segmentation is a comprehensive task to segment an object referred by a textual query from an image. In nature, the level of difficulty in this task is affected by the existence of similar objects and the complexity of the referring expression. Recent RIS models still show a significant performance gap between easy and hard scenarios. We pose that the bottleneck exists in the data, and propose a simple but powerful data augmentation method, Negative-mined Mosaic Augmentation (NeMo). This method augments a training image into a mosaic with three other negative images carefully curated by a pretrained multimodal alignment model, *e.g.*, CLIP, to make the sample more challenging. We discover that it is critical to properly adjust the difficulty level, neither too ambiguous nor too trivial. The augmented training data encourages the RIS model to recognize subtle differences and relationships between similar visual entities and to concretely understand the whole expression to locate the right target better. Our approach shows consistent improvements on various datasets and models, verified by extensive experiments.

## 1 Introduction

Referring Image Segmentation (RIS) is a fundamental task in computer vision that aims to segment objects described in a natural language expression within a given scene. Central to RIS is not merely the visual recognition of objects but also the intricate understanding of the interrelationships among these objects, interpreted through linguistic cues.

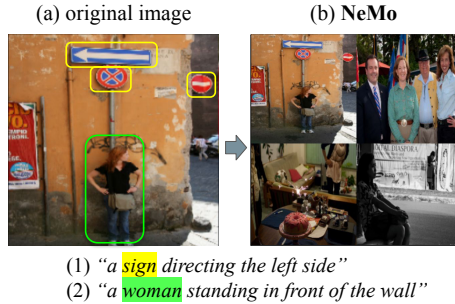
Each RIS problem often requires a different level of multimodal understanding capabilities, depending on visual ambiguity as well as linguistic complexity. For instance, having visually similar objects to the referent in an image complicates locating and identifying the correct object. In such a case, precise comprehension of the referring expression becomes a key to find the right target.

---

<sup>†</sup> Work done while at Seoul National University.

<sup>\*</sup> Equal contribution

<sup>\*\*</sup> Corresponding author



**Fig. 1:** Diverse visual and linguistic challenges of referring scenarios. In (a), query (1) demands discernment among three road signs, while query (2) involves identifying a “woman”, relatively easier due to a single instance. NeMo, our method, in (b) uses similar negative images to generate a mosaic. Query (2) becomes harder as the augmented image contains additional instances of “woman” (*e.g.*, women standing or sitting), and thus “in front of the wall” becomes crucial hint to solve the problem.

**Table 1:** Statistics of representative Referring Image Segmentation (RIS) Datasets

Dataset	RefCOCO	RefCOCO+	G-Ref
# Images	19,994	19,992	26,711
# Ref. Exp.	142,209	141,564	85,474
Query length	3.61	3.53	8.43
Obj. per query	1.76	1.67	3.03

**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. 1(a) illustrates varying degrees of difficulties even within the same image. The sentence (1) is relatively harder, since three road signs exist and the model needs to find the target by precisely understanding the entire phrase. The expression (2), on the other hand, is relatively easier, since there is only one woman, so one can easily find her regardless of the rest of the phrase.

Many existing RIS datasets, however, have not been created considering such challenge levels; rather, many examples can be solved by simply finding an object corresponding to the referred class. RefCOCO and RefCOCO+ [54], for example, contain easier scenarios with less visual ambiguity and linguistic complexity. On the contrary, G-Ref [36] is considered harder. As in Tab. 1 and Fig. 2, it contains more objects in each image on average, and queries are relatively longer.

The level of difficulty significantly varies even within the same dataset. To illustrate, we manually select 100 easy and hard samples from the G-Ref UMD test set. We select ‘easy’ samples containing only a single object per its category, straightforward to identify the target without ambiguity. For instance, in the third image in Fig. 2, there is only a single woman holding a glass, so just finding ‘a woman’ will suffice. Conversely, ‘hard’ examples contain multiple objects within the same category, necessitating a detailed perception to distinguish the intended target. Specific easy and hard examples are provided in Appendix A.

RefCOCO RefCOCO+	G-Ref (hard)	G-Ref (easy)	NeMo
			
<i>"right woman"</i> <i>"the woman"</i>	<i>"a woman getting her hair brushed"</i>	<i>"a woman holding a wine glass and is wearing a white t-shirt"</i>	

**Fig. 2:** Data samples from RIS benchmarks and augmented samples using our NeMo. RefCOCO and RefCOCO+ are characterized by relatively easier scenarios with simple referring expressions, whereas G-Refs encompass more challenging sets.

Tab. 2 shows a significant performance gap between the easy and hard samples by recent RIS models. This indicates that they are capable of picking the right object without ambiguity, but they tend to lack in understanding delicate meaning in the referring expressions and using them to distinguish multiple objects in the same (or similar) categories. To further improve the RIS performance, this observation reveals that we may need to revisit if the models have been provided with sufficiently difficult training data to learn from.

Given this problem landscape, we aim to improve the performance by tackling the data part of the training. Specifically, we postulate that amplifying the exposure of the models to challenging examples at training could fortify their capability to understand the subtle dynamics between visual and linguistic components. Such complexity often arises when multiple objects, potentially of the same class, coexist within an image, encouraging the model to fully understand both the scene and the given referring expression. However, manually labeling such ‘hard’ data examples is prohibitively expensive.

Recognizing the key factors behind the difficulty and quality of RIS training data, we introduce a simple but universally applicable data augmentation method, Negative-mined Mosaic Augmentation (NeMo). Inspired by the mosaic augmentation in YOLO v4 [3], NeMo augments each training image by combining it with three other images in a  $2 \times 2$  formation, showing four times more objects on average. However, NeMo differs from the previous method in that the extra three images are not chosen at random, but are carefully selected to create a properly challenging training example. Specifically, we propose considering relevance between the referring expression and candidate images, measured by a cross-modal retrieval model, *e.g.*, CLIP. To build a mosaic, our method selects negative images containing objects from the same or similar category to the referred object, retrieved based on the relevance to the referring expression. Augmented mosaic images mimic challenging referring examples as in Fig. 1(b), encouraging the model to learn subtle visual differences and to concretely understand the given referring expression to better locate the target.

One might concern if combining similar images may create false positives, where the correct object in an image becomes invalid due to the objects in the other quadrants. We study the possibility and impact of such false positives, and discover that it is indeed critical for the mosaic to have the right level of difficulty to be maximally effective. Based on our observations, we present a strategic retrieval process to make the mosaic neither too hard nor too easy.

From extensive experiments with five state-of-the-art RIS models, we verify that NeMo consistently improves performance across all models on multiple datasets. Furthermore, we exhibit that NeMo encourages a model to make better connections between words and visual components, grasping fine details in the scene and the referring expression. We expect our study to support the primary aim of the RIS task, distinguishing multiple candidate objects in the scene and recognizing the target based on the textual description.

Our main contributions are summarized as follows:

1. We introduce NeMo, a simple but powerful labor-free data augmentation method for Referring Image Segmentation (RIS), effective across various datasets and models.
2. We discover that it is critical to adjust the level of difficulty to successfully apply a mosaic augmentation, and propose a systematic way to tune this difficulty by generating training examples at a properly controlled difficulty.
3. We empirically verify that NeMo enhances both visual and textual understanding capabilities for segmenting the right target.

## 2 Related Work

**Referring Image Segmentation (RIS).** Existing studies [15] have concentrated on encoder and decoder architectures to handle multimodal features: RNNs [27, 29, 37, 42] and Transformers [1] for text features, and CNNs [15, 29], DeeplabV3 [1, 6, 27], and DarkNet [19, 35] for visual encoding. Transformer-based backbones [9, 19, 20, 33] and multi-scale features [5, 16, 18, 52] are recently popular to capture detailed object masks. Visual-linguistic fusion has evolved from simple concatenation [15] to syntax-based [17, 18, 53] and attention-based; to name some: LAVT [50], VPD [58], and RefSegformer [47]. VLT [10], CRIS [45], and ReSTR [23] employ cross-modal decoder fusion. ReLA [28] and CGformer [43] organizes visual features into language-conditioned tokens, capturing region-level information. PolyFormer [31] converts grounding tasks into sequential polygon generation. VPD [58] leverages a multi-scale feature map from a text-to-image diffusion model for RIS. Unlike prior studies focusing on architectural enhancements, our work redirects attention to the quality and nature of the data, proposing a straightforward augmentation to create more refined training examples.

**Datasets for RIS.** Initially, ReferIt [21], RefCOCO [54] and RefCOCO+ [54] have been introduced as benchmarks for single-target RIS. RefCOCO contains many positional words such as “front” or “the third from the right”, while RefCOCO+ prohibits such direct usage of words on positions. Many examples

in these benchmarks often present over-simplified scenarios with a short query, without enough ambiguities in images. In contrast, recent datasets embrace more complex scenes and intricate linguistic expressions. G-Ref [36] contains relatively longer sentences, making the task more challenging. Built upon [24], Phrase-Cut [46] provides examples with multiple targets with attributes, categories, and relationships in phrases. GRES [28] steps up the complexity by incorporating references to no target and multiple objects. CGFormer [43] introduces new splits on RefCOCO datasets to measure generalization capability on objects of unseen classes. Our method is aligned with these studies for better data quality, but we directly generate complex examples via augmentation.

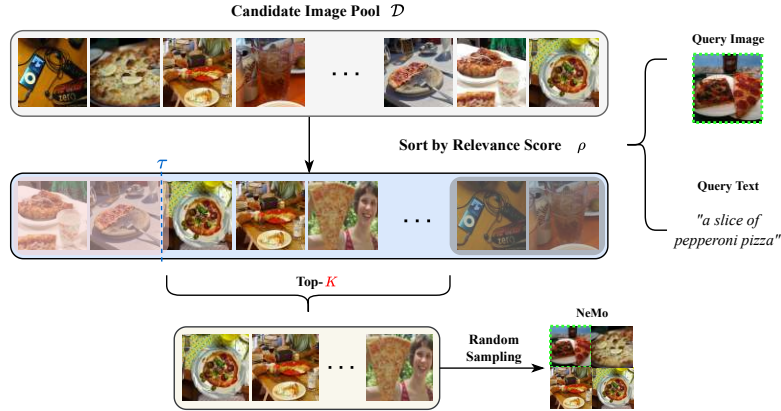
**Data Augmentation Methods.** Various data augmentation methods have been proposed for semantic segmentation. Early strategies involve random object pasting [11]. Multi-modal mixup [13, 57] fuses images and associated texts. Beyond this, CutMix [56] transposes rectangular image sections onto others. Mosaic data augmentation, originally pioneered for object detection by YOLO v4 [3], generates a composite image from segments of multiple sources, preserving their ground truths. MixGen [13] is another multimodal augmentation, blending two images and combining pairs of text sequences. However, this mosaic-based method has not been specifically designed for the RIS task, to the best of our knowledge. In this paper, we discover that simply putting four images into a mosaic does not guarantee maximal performance improvement, and propose a way specific to RIS to curate mosaics.

**Retrieval Augmented Vision-Language Models.** Retrieval Augmented models [7] employ retrieval to leverage additional knowledge from external data to enhance model’s learning capabilities. Starting from KNN-LM [22], several works in natural language processing have enhanced large language models by connecting them with external sources, intricately arrayed in structured syntax and semantic relations [4, 12, 25, 32, 39, 55]. This approach has expanded to various vision-language tasks; *e.g.*, image synthesis [2, 8, 41], classification [34], and multi-modal applications [51]. This method is used to generate hard negatives in tasks like multimedia event extraction [26] and provide task-specific data augmentation [30, 34]. Additionally, RA-CLIP [48], K-Lite [40], and ASIF [38] employ retrieval to enrich visual concepts and to align image-text modalities.

In our work, we employ retrieval for mosaic augmentation tailored for the RIS task. By generating ambiguous scenarios with hard negatives, we enforce the model to better connect textual expressions with visual components, aiming to improve the model performance in complex situations.

### 3 NeMo: Negative-mined Mosaic Augmentation

In this section, we propose a model-agnostic data augmentation framework, Negative-mined Mosaic Augmentation (NeMo), for the RIS task. Particularly, we introduce a simple but flexible augmentation technique that guides us in choosing proper negative images to create harder training samples. Fig. 3 illustrates the overall pipeline of our approach.



**Fig. 3: Overall NeMo pipeline.** Given an image and a query, it selects negative images at a proper level of difficulty, filtering out visually or semantically images to the query to avoid false negatives and irrelevant (easy) images identified by text-to-image retrieval. It randomly selects three among the remaining to construct a mosaic.

### 3.1 Motivation for Negative Mining

As discussed in Sec. 1, the level of ambiguity in images largely depends on whether they contain objects of similar categories related to the referring expression. We leverage this nature of the RIS task to control the difficulty of each problem instance by carefully mining the negatives in the mosaic. Fig. 4 demonstrates how the task difficulty varies with the visual and contextual relationships between the negative images and the text query. For example, a distinct positive image of pizza combined with unrelated images like flowers or buses in (a) poses little challenge for the model. Conversely, images of multiple skateboarders jumping in (c) provide too many plausible choices for the given phrase, leading to confusion. We aim to craft training samples with the right level of difficulty as in (b), neither too easy nor too hard. Balancing challenge and distinguishability can foster fine-grained learning of the relationships between the referring expression and the objects in the scene.

To select suitable negative images, we define two hyperparameters as illustrated in Fig. 3: uni-modal or cross-modal similarity score threshold ( $\tau$ ) to exclude overly similar examples, and the number of top negative image candidates ( $K$ ) to consider. The filtering process is detailed below.

### 3.2 Negative Image Mining Methodology

For an original training example  $(I, T)$ , where  $I$  and  $T$  stand for the original image and referring expression, our approach aims to retrieve negative images that are moderately aligned with  $T$  from the pool of all images,  $\mathcal{D}$ . To quantify the relevance of each candidate image  $I^{(i)} \in \mathcal{D}$ , we rely on visual-text similarity scores estimated by a pre-trained cross-modal model, *e.g.*, CLIP [26].



**Fig. 4: Comparison of negative image choices** Finding the “rightmost pizza” in (a) is nearly as *easy* as in the single image, as there is no other pizza-like object. Multiple road signs in (b) require discerning the relative location of a woman and an SUV, more *challenging* than the original single image. (c) is *invalid* as multiple images contain “a man jumping with a skateboard”.

**Determining the Upper-bound.** At a glance, it looks straightforward to choose the top  $N$  images that are most relevant to the target text  $T$  as hard negatives. The images with the highest relevance in  $\mathcal{D}$ , however, can be visually too similar to the query, which may result abundant potential choices. Those negatives are in fact ‘false negatives’, where there exists another object perfectly in accordance with the referred expression  $T$  within the retrieved negative image and thus it becomes impossible to find the intended one, as in Fig. 4(c).

To address this, we filter out potential false negatives, or excessively relevant candidate images to  $T$ . Specifically, we compute the relevance score  $\rho^{(i)}$  of each candidate image  $I^{(i)} \in \mathcal{D}$  with the referring expression  $T$  by  $\mathbf{t}^\top \mathbf{v}^{(i)}$ , where  $\mathbf{v}^{(i)}$  and  $\mathbf{t}$  are CLIP visual and text embedding of the candidate image  $I^{(i)}$  and  $T$ , respectively. NeMo then prevents potential false negatives by excluding candidate images that are too relevant to  $T$  with  $\rho^{(i)} \geq \tau$ , where  $\tau$  is a hyperparameter to control the tolerance of upper-bound filtering. With a large  $\tau$ , we filter only extremely relevant images from the candidate set, while with a smaller  $\tau$ , we aim to filter more to ensure less false negatives to occur.

Alternatively, we may use  $\rho$  for the image-to-image similarities,  $\mathbf{v}^\top \mathbf{v}^{(i)}$ , between the positive image  $I$  and all other negative candidates  $I^{(i)} \in \mathcal{D}$ , where  $\mathbf{v}$  is the CLIP visual embedding of  $I$ . This can be useful to capture a highly relevant candidate image captioned with a less similar text form. It is also possible to use both text-to-image (t-i) and image-to-image (i-i) similarities, filtering out images either  $\rho_{\text{t-i}}^{(i)} \geq \tau_{\text{t-i}}$  or  $\rho_{\text{i-i}}^{(i)} \geq \tau_{\text{i-i}}$ . See Appendix H for empirical comparison.

**Determining the Lower-bound.** After we filter out excessively similar images that are  $\rho^{(i)} \geq \tau$ , we collect  $K$  most plausible images described by  $T$  among the remaining candidate images. At this step, the relevance of an image  $I^{(i)}$  can be computed using the same  $\rho$  or some other  $\rho'$ . To keep the framework general, we

use  $\rho'$  for the relevance used in this step; that is, we collect the top  $K$  images with highest  $\rho'$  from the set  $\{I^{(i)} \in \mathcal{D} : \rho^{(i)} < \tau\}$ . We illustrate the simpler case with  $\rho = \rho'$  in Fig. 3, and the general case with  $\rho \neq \rho'$  in Appendix H. We then randomly select 3 out of the  $K$  candidates, arrange them with the positive  $I$  in a  $2 \times 2$  mosaic grid, and resize the resulting image to half width and height. The quadrant corresponding to the positive image is labeled accordingly, while the other three quadrants are set to negative for all pixels.

$K$  is another hyperparameter to control the difficulty of an augmented image. Lower  $K$  would select more plausible images related to  $T$ , which might include other partially correct objects. With a higher  $K$ , chances of choosing less relevant images increase. When  $K \approx |\mathcal{D}|$ , it is equivalent to the uniform selection.

**Augmentation Ratio  $\gamma$ .** In practice, we expose the model to single images as well as augmented ones during training, since eventually the model performs on single images. Specifically, we apply NeMo with a probability of  $\gamma \in [0, 1]$ , while the rest  $(1 - \gamma)$  uses an original single image.

**Summary.** The overall process ensures a proper level of difficulty, guided by the parameters  $\tau$  and  $K$ . When adjusted properly,  $\tau$  helps to filter out images that are too similar, and  $K$  determines the number of relevant images to consider. Careful calibration of them guides the chosen images to be sufficiently similar to  $T$  but at the same time distinct enough from the positive image  $I$ . Such a balanced selection of images helps the model to better understand and adapt to various visual contexts, thereby significantly improving its learning capabilities.

### 3.3 Addressing False Negatives & False Positives

Even with a careful choice of  $\tau$  and  $K$ , NeMo may still generate false negatives (FN) and false positives (FP), especially on simpler referring expressions.

First, FN occurs when another object perfectly matching with  $T$  exists in one of the chosen negative images. This seems problematic since it is still labeled negative but it is essentially a positive. Nevertheless, we claim that FNs are not significantly detrimental to performance. Similar to the Masked Language Modeling in BERT [9], where multiple valid fill-ins exist for a blank but only one is taught, the model would learn the probability distribution of plausible objects upon sufficient examples and repeated training.

Meanwhile, FP happens when the correct object in  $I$  is affected by its placement in relation to other images, often due to positional descriptors<sup>5</sup> in  $T$ . This is more common when the target object is designated by its relationship with other objects or its position within the image frame. For example, if “a woman in the left” is positioned in the right quadrant within the  $2 \times 2$  grid and another woman appears in the left negative image, the target designation shifts, invalidating the label on the original target.

FPs can be a more critical problem than FNs, since they can *mislead* the model to select a wrong object. We show that our NeMo inherently provides a

<sup>5</sup> e.g., left, right, low, high, top, bottom, o'clock, corner

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

RIS model	NeMo	RefCOCO (UNC)			RefCOCO+ (UNC)			G-Ref (UMD)		GRES	Average Gain
		Val	TestA	TestB	Val	TestA	TestB	Val	Test	Val	
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		$+1.11 \pm 0.79$			$+1.21 \pm 0.48$			$+1.55 \pm 0.99$		$+3.11 \pm 2.83$	

way to prevent FPs; choosing a lower  $\tau$  would reduce the chance of FP by filtering out even less relevant candidate images. See Sec. 4.4 for empirical verification.

## 4 Experiments

### 4.1 Experimental Settings

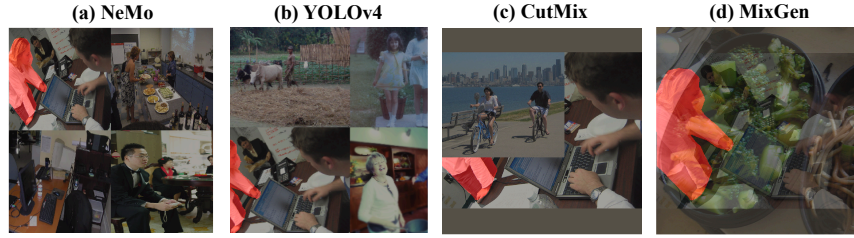
**Datasets.** We evaluate on four widely-used RIS benchmarks: RefCOCO [21], RefCOCO+ [21], G-Ref [36], and GRES [28]. Dataset statistics are summarized in Tab. 1. For RefCOCO and RefCOCO+, we use the train, validation, test A and B partitions. We use both UMD [54] and Google [36] split for G-Ref dataset. Refer to Appendix D for the Google split.

**Metrics.** We adopt three metrics. First, overall intersection over union is the proportion of the intersection area to the union area across all test samples. Due to its tendency to favor larger objects, we also use mean intersection over union, representing the average intersection between the prediction and the ground truth for all samples. Lastly, we report Precision@ $p$ , the ratio of samples whose IoU with the ground truth exceeds the threshold  $p$ , with  $p \in \{0.5, 0.7, 0.9\}$ .

**RIS Models.** To verify the applicability of our data augmentation method, we experiment with five state-of-the-art RIS models: LAVT [50], CRIS [45], ReLA [28], CGFormer [43] and VPD [58]. For each method, we compare the overall RIS performance with and without applying our NeMo. Refer to Appendix B for implementation details.

### 4.2 Effectiveness of the Proposed Method

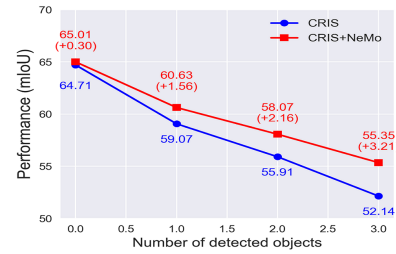
Tab. 3 compares the performance of various RIS models with and without applying our method in Overall IoU (see Appendix C for Mean IoU). We observe that NeMo improves the performance under all settings, across all datasets regardless of the RIS model. This indeed reveals that a bottleneck has also been in



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

**Fig. 5: Samples from other multi-modal augmentation method.**  $2 \times 2$  mosaic samples from (a)NeMo, (b)YOLOv4, (c)CutMix, (d)MixGen.**Table 4: Comparison to other multi-modal augmentations on G-Ref**

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>

**Fig. 6: mIoUs for the different number of negative objects.**

data, not just the modeling aspect, and our method effectively provides ambient training examples with a curated complexity.

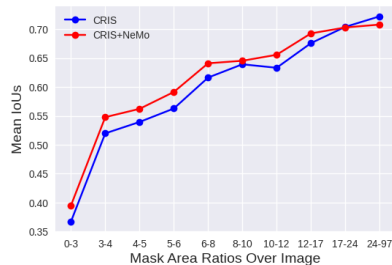
An interesting observation is that the degree of improvement differs by the datasets. On average, we observe a larger performance boost on a more complex dataset (G-Ref) than simpler ones (RefCOCO and RefCOCO+). Compared to the other two, G-Ref contains nearly twice many objects and three times longer queries (see Tab. 1). G-Ref benefits more from NeMo because of its intricate referring expressions and visually dense scenes. Similarly, on GRES, extension of RefCOCO+ with no-target and multi-target queries, NeMo shows significant performance gain. The complex expressions in GRES amplifies the challenge of correctly identifying target(s), and the result validates the robustness of our method in handling more complex referring scenarios. In contrast, our approach provides less impact on the relatively simpler datasets, because they do not require a nuanced understanding to differentiate similar objects. Nevertheless, we still notice some improvements even on these simpler datasets, indicating that NeMo is still beneficial and does not affect negatively.

### 4.3 Comparison with Other Augmentation Methods

Tab. 4 compares different multi-modal augmentation methods on G-Ref. We observe that most methods degrade performance, showing that other augmentation

**Table 5:** Performance over various sentence lengths on G-Ref UMD test split.

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

**Fig. 7: mIoUs for various object sizes.** Sizes are binned such that each bin contains 10% of the test set.

methods are not suitable for the RIS task. Fig. 5 illustrates how they can fail for RIS; YOLOv4 and CutMix often lose or obstruct referents after cropping or overlaying the image. MixGen also underperforms, likely due to difficulties understanding the whole scene of the original image while interpolated. These results endorse NeMo as a suitable augmentation approach for the RIS task.

#### 4.4 Detailed Analysis of the Proposed Method

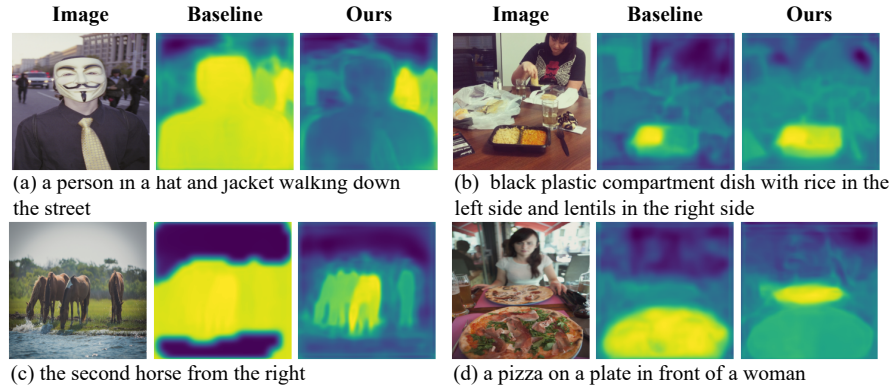
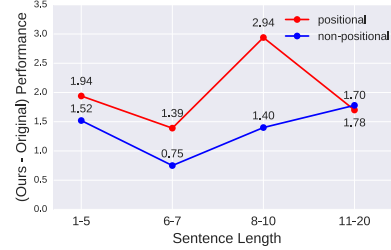
**Performance on Challenging Scenarios.** RIS task is challenging when it cannot be easily resolved through class or positional keywords alone. This occurs when there exist multiple objects of the same class as the referent, where we expect NeMo to be particularly effective. To evaluate this systematically, we compare the performance with varying number of negative objects within the image. We first detect objects in the scene with a pre-trained detector, and identify same-class objects that closely overlap with the target to count negative objects. Fig. 6 shows that the performance gap between with and without NeMo gets larger with more negative objects in the image. This indicates that NeMo performs better on challenging samples as our expectation.

**Robustness on Object Scale.** We evaluate the impact of our method across various object sizes in Fig. 7. Improvement is observed in most cases, especially for smaller objects. This can be attributed to the wider range of object scales seen during training by integrating objects both in the original and 1/4 size. Fig. 9(a) demonstrates that our method successfully detects a blurred and small “person” in the background, positioned behind the most prominent person in the image. Overall, we observe clearer boundaries in the final activation maps with NeMo, commonly illustrated in Fig. 9.

**Complexity of Referring Expression.** Following [23, 29], we measure how our method behaves depending on sentence lengths. Tab. 5 shows that NeMo is effective across all sentence lengths, even with longer complex ones. NeMo also helps capture important linguistic cues for grounding. In Fig. 9(b), our method

**Fig. 8:** Performance gain for queries with and without positional keywords.**Table 6:** Effect of FP on RefCOCO

Method	Val	oIoU		Prec (Val)	
		TestA	TestB	0.5	0.7
CRIS	66.68	70.62	59.93	80.89	69.23
Ours	<b>68.66</b>	<b>72.82</b>	<b>63.33</b>	<b>82.87</b>	<b>70.63</b>
Ours wo/ FP	68.30	72.41	63.26	81.96	70.18

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

segments the entire dish while the baseline only detects the left half without fully understanding the query describing the right half.

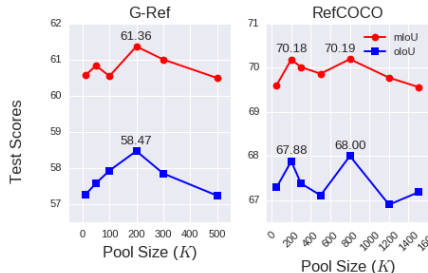
**Positional Understanding.** Fig. 8 demonstrates that NeMo exhibits stronger improvements especially for queries with positional keywords even when the queries are long and complex. Along with the improved visual and textual understanding with NeMo, the model is exposed to various samples in mosaic not confined to the specific local positions, which requires deeper positional understanding. As seen in Fig. 9(b,d), our method captures objects more accurately in scenarios involving directional expressions, indicating improved understanding both in absolute and relative positions. In Fig. 9(c), it is apparent that our method yields an output with more distinct shape for “the second” horse.

**Effect of False Positives.** Even after tuning  $\tau$ , one might concern non-zero probability of FPs. To address this, we conduct a rule-based experiment to restrict the location of the positive image, if the text contains positional keywords. For example, images labeled with ‘left’ are forced to be in one of the left quadrants. This experiment is carried out on RefCOCO, featured by simpler texts. As shown in Tab. 6, performance gap between with and without the constraint

**Fig. 10: Ablation on  $K$** 

**Table 7: Ablation on  $\gamma$**

$\gamma$	P@0.5	P@0.7	P@0.9	oIoU	mIoU
0.8	69.61	56.62	14.64	57.31	60.62
0.6	69.00	57.09	16.06	<b>58.15</b>	<b>60.90</b>
0.5	69.51	<b>57.23</b>	<b>16.38</b>	57.81	60.85
0.4	69.38	56.33	15.85	57.40	60.88
0.2	<b>70.16</b>	56.46	15.38	57.38	60.62

**Table 8: Ablation on  $\tau$** 

$\tau$	G-Ref(Val) $K = 200$					RefCOCO(Val) $K = 800$				
	P@0.5	P@0.7	P@0.9	oIoU	mIoU	P@0.5	P@0.7	P@0.9	oIoU	mIoU
1.00	69.14	55.76	15.09	58.06	60.73	81.90	69.36	19.96	<b>68.04</b>	70.34
0.85	69.47	56.29	15.08	57.57	60.50	<b>82.14</b>	<b>71.04</b>	<b>20.07</b>	68.01	<b>70.71</b>
0.75	<b>70.01</b>	<b>56.60</b>	<b>15.89</b>	<b>58.47</b>	<b>61.36</b>	81.35	70.23	19.81	68.00	70.19
0.60	69.63	56.17	14.07	57.92	60.36	81.66	70.02	19.31	67.74	70.29

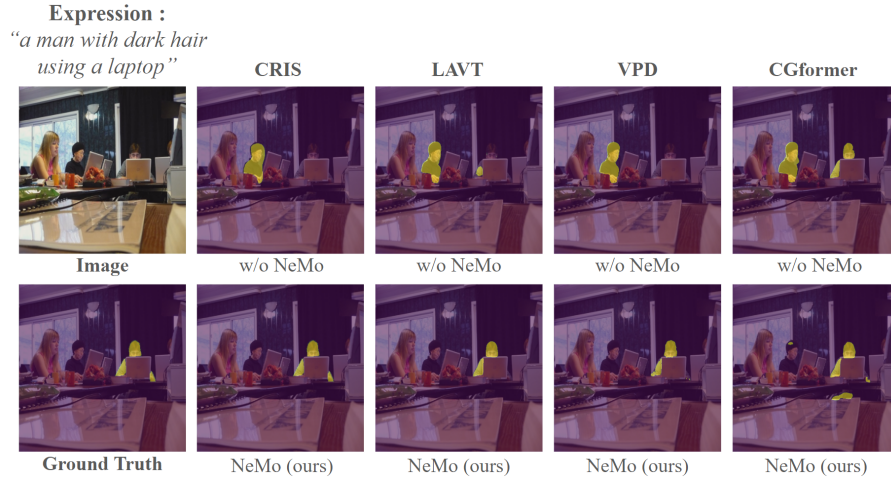
turns out to be minimal, indicating that  $\tau$  and  $K$  is effective in mitigating FPs. Intriguingly, NeMo without the constraint performs the best. We interpret this as 1) NeMo hardly creates FPs with well-tuned  $\tau$  and  $K$ , and 2) even if some FPs exist, the benefit from allowing more diversity of augmented images is bigger. See more discussion in Appendix J.

#### 4.5 Ablation Study

We conduct an ablation study to break down the contributions of individual components of NeMo. All experiments are conducted with CRIS [45] on the UMD validation set of G-Ref and UNC validation set of RefCOCO.

**Augmentation Ratio  $\gamma$ .** Recall from Sec. 3.2 that we apply the proposed augmentation with the probability of  $\gamma$  to each training sample. When NeMo is not applied, original single images are used for training as the model eventually infers on the single images. Before adjusting the difficulty of mosaics, we conduct this experiment at the median level of difficulty by setting  $K = |\mathcal{D}|/2$ . Tab. 7 indicates that best performance is obtained when  $\gamma$  is between 0.5 to 0.6, while performance slightly drops with  $\gamma$  too large or too small.

**Negative Image Pool Size  $K$ .** We explore the optimal pool size of the candidate negative images,  $K$ , on both G-Ref and RefCOCO. Adjusting  $K$  tunes not just the number, but also the minimum relevance of the negative images to be considered. The other hyperparameters are set to the default values,  $\tau = 0.75$  and  $\gamma = 0.6$ . Fig. 10 reveals that both exceedingly low and high  $K$  are suboptimal. With too small  $K$ , only mostly similar images may remain, creating FP and FN frequently. In contrast, an extremely large  $K$  brings little additional challenge, getting closer to the uniform selection. This confirms that a moderate level



**Fig. 11:** Visualization of predictions from the G-Ref test set.

of challenge in the negative samples is most beneficial. Besides, the optimal  $K$  values differ significantly by the datasets: 200 for G-Ref and 800 for RefCOCO. In a simpler dataset, allowing less visually similar images would enhance overall performance by minimizing the chances to cause FP/FN.

**Upper-bound Threshold  $\tau$ .** We explore the optimal  $\tau \in \{0.6, 0.75, 0.85, 1.0\}$  where  $K$  is fixed to the best value found for each dataset above. Tab. 8 shows that a moderate  $\tau$  around 0.75 to 0.85 yields the best performance. With a larger  $\tau$ , less images are filtered out, causing more chances for FP/FN. With too small  $\tau$ , it becomes closer to random sampling again. We further explore the cross-effect of  $\tau$  and  $K$  and mosaic design choices in Appendix G.

## 5 Conclusion

This paper proposes NeMo, an advanced mosaic augmentation method that exposes an RIS model to more intricate and challenging examples, thereby enhancing its visual and linguistic understanding for locating and segmenting the referent. NeMo brings consistent performance improvement over various state-of-the-art RIS models on multiple datasets, especially on datasets with higher visual-linguistic complexity. Although NeMo shows consistent improvement, data augmentation in this context still remains a relatively unexplored territory. More sophisticated methodologies, such as an object-level parsing, could potentially further enhance retrieval, but we leave this for future research.

**Limitations.** Our method performs well when the dataset contains relatively homogeneous images. If it contains images from diverse domains (*e.g.*, X-ray, sketches, or satellite images), mosaic creation may result in unnatural combinations, potentially degrading the performance.

## Acknowledgements

We appreciate Jeongwoo Shin for insightful discussion. This work was supported by the New Faculty Startup Fund from Seoul National University, by Samsung Electronics Co., Ltd (IO230414-05943-01, RAJ0123ZZ-80SD), by Youlchon Foundation (Nongshim Corp.), and by National Research Foundation (NRF) grants (No. 2021H1D3A2A03038607/50%, RS-2024-00336576/10%, RS-2023-00222663/5%) and Institute for Information & communication Technology Planning & evaluation (IITP) grants (No. RS-2024-00353131/25%, RS-2022-II220264/10%), funded by the government of Korea.

## References

1. Bellver, M., Ventura, C., Silberer, C., Kazakos, I., Torres, J., Giro-i Nieto, X.: A closer look at referring expressions for video object segmentation. *Multimedia Tools and Applications* **82**(3), 4419–4438 (2023) [4](#)
2. Blattmann, A., Rombach, R., Oktay, K., Müller, J., Ommer, B.: Retrieval-augmented diffusion models. In: *NeurIPS* (2022) [5](#)
3. Bochkovskiy, A., Wang, C.Y., Liao, H.Y.M.: YOLOv4: Optimal speed and accuracy of object detection. *arXiv:2004.10934* (2020) [3](#), [5](#), [10](#), [vii](#)
4. Borgeaud, S., Mensch, A., Hoffmann, J., Cai, T., Rutherford, E., Millican, K., Van Den Driessche, G.B., Lespiau, J.B., Damoc, B., Clark, A., et al.: Improving language models by retrieving from trillions of tokens. In: *ICML* (2022) [5](#)
5. Chen, D.J., Jia, S., Lo, Y.C., Chen, H.T., Liu, T.L.: See-through-text grouping for referring image segmentation. In: *ICCV* (2019) [4](#)
6. Chen, L.C., Papandreou, G., Schroff, F., Adam, H.: Rethinking atrous convolution for semantic image segmentation. *arXiv:1706.05587* (2017) [4](#)
7. Chen, W., Hu, H., Chen, X., Verga, P., Cohen, W.W.: MuRAG: Multimodal retrieval-augmented generator for open question answering over images and text. In: *EMNLP* (2022) [5](#)
8. Chen, W., Hu, H., Saharia, C., Cohen, W.W.: Re-Imagen: Retrieval-augmented text-to-image generator. *arXiv:2209.14491* (2022) [5](#)
9. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: BERT: Pre-training of deep bidirectional transformers for language understanding. In: *NAACL* (2019) [4](#), [8](#), [i](#)
10. Ding, H., Liu, C., Wang, S., Jiang, X.: Vision-language transformer and query generation for referring segmentation. In: *ICCV* (2021) [4](#)
11. Ghiasi, G., Cui, Y., Srinivas, A., Qian, R., Lin, T.Y., Cubuk, E.D., Le, Q.V., Zoph, B.: Simple copy-paste is a strong data augmentation method for instance segmentation. In: *CVPR* (2021) [5](#)
12. Guu, K., Lee, K., Tung, Z., Pasupat, P., Chang, M.: Retrieval augmented language model pre-training. In: *ICML* (2020) [5](#)
13. Hao, X., Zhu, Y., Appalaraju, S., Zhang, A., Zhang, W., Li, B., Li, M.: MixGen: A new multi-modal data augmentation. In: *WACV* (2023) [5](#), [10](#)
14. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *CVPR* (2016) [i](#)
15. Hu, R., Rohrbach, M., Darrell, T.: Segmentation from natural language expressions. In: *ECCV* (2016) [4](#)
16. Hu, Z., Feng, G., Sun, J., Zhang, L., Lu, H.: Bi-directional relationship inferring network for referring image segmentation. In: *CVPR* (2020) [4](#)

17. Huang, S., Hui, T., Liu, S., Li, G., Wei, Y., Han, J., Liu, L., Li, B.: Referring image segmentation via cross-modal progressive comprehension. In: CVPR (2020) [4](#)
18. Hui, T., Liu, S., Huang, S., Li, G., Yu, S., Zhang, F., Han, J.: Linguistic structure guided context modeling for referring image segmentation. In: ECCV (2020) [4](#)
19. Jing, Y., Kong, T., Wang, W., Wang, L., Li, L., Tan, T.: Locate then segment: A strong pipeline for referring image segmentation. In: CVPR (2021) [4](#)
20. Kamath, A., Singh, M., LeCun, Y., Synnaeve, G., Misra, I., Carion, N.: Mdetr-modulated detection for end-to-end multi-modal understanding. In: ICCV (2021) [4](#)
21. Kazemzadeh, S., Ordonez, V., Matten, M., Berg, T.: ReferItGame: Referring to objects in photographs of natural scenes. In: EMNLP (2014) [4, 9](#)
22. Khandelwal, U., Levy, O., Jurafsky, D., Zettlemoyer, L., Lewis, M.: Generalization through memorization: Nearest neighbor language models. arXiv:1911.00172 (2019) [5](#)
23. Kim, N., Kim, D., Lan, C., Zeng, W., Kwak, S.: Restr: Convolution-free referring image segmentation using transformers. In: CVPR (2022) [4, 11](#)
24. Krishna, R., Zhu, Y., Groth, O., Johnson, J., Hata, K., Kravitz, J., Chen, S., Kalantidis, Y., Li, L.J., Shamma, D.A., et al.: Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision* **123**, 32–73 (2017) [5](#)
25. Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W.t., Rocktäschel, T., et al.: Retrieval-augmented generation for knowledge-intensive NLP tasks. *NeurIPS* (2020) [5](#)
26. Li, M., Xu, R., Wang, S., Zhou, L., Lin, X., Zhu, C., Zeng, M., Ji, H., Chang, S.F.: Clip-event: Connecting text and images with event structures. In: CVPR (2022) [5, 6](#)
27. Li, R., Li, K., Kuo, Y.C., Shu, M., Qi, X., Shen, X., Jia, J.: Referring image segmentation via recurrent refinement networks. In: CVPR (2018) [4](#)
28. Liu, C., Ding, H., Jiang, X.: GRES: Generalized referring expression segmentation. In: CVPR (2023) [4, 5, 9, 11, iv](#)
29. Liu, C., Lin, Z., Shen, X., Yang, J., Lu, X., Yuille, A.: Recurrent multimodal interaction for referring image segmentation. In: ICCV (2017) [4, 11](#)
30. Liu, H., Son, K., Yang, J., Liu, C., Gao, J., Lee, Y.J., Li, C.: Learning customized visual models with retrieval-augmented knowledge. In: CVPR (2023) [5](#)
31. Liu, J., Ding, H., Cai, Z., Zhang, Y., Satzoda, R.K., Mahadevan, V., Manmatha, R.: PolyFormer: Referring image segmentation as sequential polygon generation. In: CVPR (2023) [4](#)
32. Liu, W., Zhou, P., Zhao, Z., Wang, Z., Ju, Q., Deng, H., Wang, P.: K-BERT: Enabling language representation with knowledge graph. In: AAAI (2020) [5](#)
33. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. In: ICCV (2021) [4, i](#)
34. Long, A., Yin, W., Ajanthan, T., Nguyen, V., Purkait, P., Garg, R., Blair, A., Shen, C., van den Hengel, A.: Retrieval augmented classification for long-tail visual recognition. In: CVPR (2022) [5](#)
35. Luo, G., Zhou, Y., Sun, X., Cao, L., Wu, C., Deng, C., Ji, R.: Multi-task collaborative network for joint referring expression comprehension and segmentation. In: CVPR (2020) [4](#)
36. Mao, J., Huang, J., Toshev, A., Camburu, O., Yuille, A.L., Murphy, K.: Generation and comprehension of unambiguous object descriptions. In: CVPR (2016) [2, 5, 9](#)

37. Margffoy-Tuay, E., Pérez, J.C., Botero, E., Arbeláez, P.: Dynamic multimodal instance segmentation guided by natural language queries. In: ECCV (2018) 4
38. Norelli, A., Fumero, M., Maiorca, V., Moschella, L., Rodola, E., Locatello, F.: ASIF: Coupled data turns unimodal models to multimodal without training. In: NeurIPS (2023) 5
39. Peters, M.E., Neumann, M., Logan IV, R.L., Schwartz, R., Joshi, V., Singh, S., Smith, N.A.: Knowledge enhanced contextual word representations. In: EMNLP (2019) 5
40. Shen, S., Li, C., Hu, X., Xie, Y., Yang, J., Zhang, P., Gan, Z., Wang, L., Yuan, L., Liu, C., et al.: K-LITE: Learning transferable visual models with external knowledge. In: NeurIPS (2022) 5
41. Sheynin, S., Ashual, O., Polyak, A., Singer, U., Gafni, O., Nachmani, E., Taigman, Y.: KNN-diffusion: Image generation via large-scale retrieval. In: ICLR (2023) 5
42. Shi, H., Li, H., Meng, F., Wu, Q.: Key-word-aware network for referring expression image segmentation. In: ECCV (2018) 4
43. Tang, J., Zheng, G., Shi, C., Yang, S.: Contrastive grouping with transformer for referring image segmentation. In: CVPR (2023) 2, 4, 5, 9, 11, iv
44. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. In: NIPS (2017) i
45. Wang, Z., Lu, Y., Li, Q., Tao, X., Guo, Y., Gong, M., Liu, T.: CRIS: Clip-driven referring image segmentation. In: CVPR (2022) 2, 4, 9, 11, 13, iv, x
46. Wu, C., Lin, Z., Cohen, S., Bui, T., Maji, S.: PhraseCut: Language-based image segmentation in the wild. In: CVPR (2020) 5
47. Wu, J., Li, X., Li, X., Ding, H., Tong, Y., Tao, D.: Towards robust referring image segmentation. IEEE Transactions on Image Processing (2024) 4
48. Xie, C.W., Sun, S., Xiong, X., Zheng, Y., Zhao, D., Zhou, J.: RA-CLIP: Retrieval augmented contrastive language-image pre-training. In: CVPR (2023) 5
49. Yang, S., Xia, M., Li, G., Zhou, H.Y., Yu, Y.: Bottom-up shift and reasoning for referring image segmentation. In: CVPR (2021) i
50. Yang, Z., Wang, J., Tang, Y., Chen, K., Zhao, H., Torr, P.H.: LAVT: Language-aware vision transformer for referring image segmentation. In: CVPR (2022) 2, 4, 9, 11, iv
51. Yasunaga, M., Aghajanyan, A., Shi, W., James, R., Leskovec, J., Liang, P., Lewis, M., Zettlemoyer, L., Yih, W.t.: Retrieval-augmented multimodal language modeling. arXiv:2211.12561 (2023) 5
52. Ye, L., Rochan, M., Liu, Z., Wang, Y.: Cross-modal self-attention network for referring image segmentation. In: CVPR (2019) 4
53. Yu, L., Lin, Z., Shen, X., Yang, J., Lu, X., Bansal, M., Berg, T.L.: MAttNet: Modular attention network for referring expression comprehension. In: CVPR (2018) 4
54. Yu, L., Poirson, P., Yang, S., Berg, A.C., Berg, T.L.: Modeling context in referring expressions. In: ECCV (2016) 2, 4, 9
55. Yu, W., Zhu, C., Fang, Y., Yu, D., Wang, S., Xu, Y., Zeng, M., Jiang, M.: DictBERT: Enhancing language model pre-training with dictionary. In: ACL (2021) 5
56. Yun, S., Han, D., Oh, S.J., Chun, S., Choe, J., Yoo, Y.: CutMix: Regularization strategy to train strong classifiers with localizable features. In: ICCV (2019) 5, 10
57. Zhang, H., Cisse, M., Dauphin, Y.N., Lopez-Paz, D.: mixup: Beyond empirical risk minimization. arXiv:1710.09412 (2017) 5
58. Zhao, W., Rao, Y., Liu, Z., Liu, B., Zhou, J., Lu, J.: Unleashing text-to-image diffusion models for visual perception. arXiv:2303.02153 (2023) 4, 9, 11, iv