

Unveiling Parts Beyond Objects: Towards Finer-Granularity Referring Expression Segmentation

Wenxuan Wang^{1,2,3*} Tongtian Yue^{1,2*} Yisi Zhang⁴ Longteng Guo¹ Xingjian He¹
Xinlong Wang³ Jing Liu^{1,2†}

¹ Institute of Automation, Chinese Academy of Sciences (CASIA)

² School of Artificial Intelligence, University of Chinese Academy of Sciences (UCAS)

³ Beijing Academy of Artificial Intelligence (BAAI)

⁴ University of Science and Technology Beijing (USTB)

{wangwenxuan2023@ia.ac.cn, wangxinlong@baai.ac.cn, jliu@nlpr.ia.ac.cn}

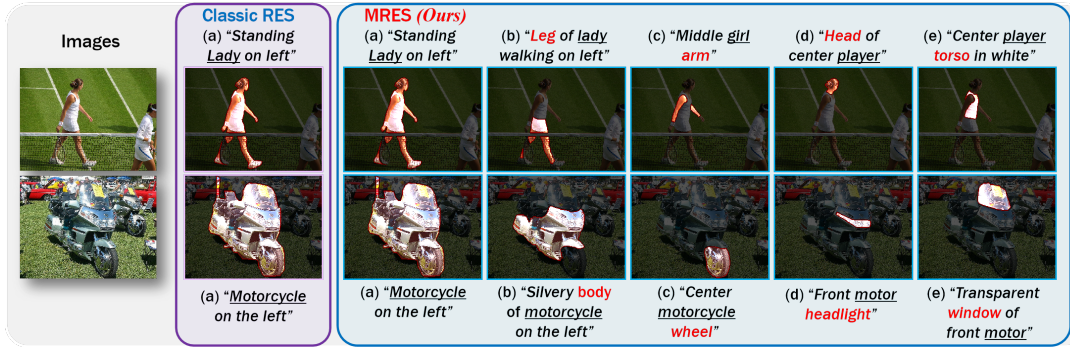


Figure 1. Classic Referring Expression Segmentation (RES) only supports expressions that indicate a single target object, e.g., (a). Compared with classic RES, the proposed **Multi-Granularity Referring Expression Segmentation (MRES)** task supports expressions indicating the specific *part-level regions* of target objects, e.g., part-level expressions like (b)-(e) from our newly built RefCOCO_m benchmark.

Abstract

Referring expression segmentation (RES) aims at segmenting the foreground masks of the entities that match the descriptive natural language expression. Previous datasets and methods for classic RES task heavily rely on the prior assumption that one expression must refer to object-level targets. In this paper, we take a step further to finer-grained part-level RES task. To promote the object-level RES task towards finer-grained vision-language understanding, we put forward a new multi-granularity referring expression segmentation (MRES) task and construct an evaluation benchmark called RefCOCO_m by manual annotations. By employing our automatic model-assisted data engine, we build the largest visual grounding dataset namely MRES-32M, which comprises over 32.2M high-quality masks and captions on the provided 1M images. Besides, a simple yet strong model named UniRES is designed to accomplish the unified object-level and part-level grounding task. Extensive experiments on our RefCOCO_m for MRES and three datasets (i.e., RefCOCO(+g)) for classic RES task demon-

strate the superiority of our method over previous state-of-the-art methods. To foster future research into fine-grained visual grounding, our benchmark RefCOCO_m, the MRES-32M dataset and model UniRES will be publicly available at <https://github.com/Rubics-Xuan/MRES>.

1. Introduction

As one of the most challenging tasks in vision-language understanding, referring expression segmentation (RES) aims to locate specific regions at the pixel level based on a descriptive language expression. Compared to traditional visual segmentation tasks that focus on images or videos alone, RES poses greater difficulties and challenges due to the necessity of strong comprehension across modalities, but it can simultaneously alleviate the problem of pre-defined categories in conventional object detection or segmentation. With the real-world scene often requiring diversity in target identification, RES task holds vast potential for applications, e.g., language-based human-object interaction and interactive image editing.

Since the concept of RES task was initially proposed in

*Equal contribution.

†Corresponding author.

Table 1. Comparison among different referring expression datasets, including ReferIt [18], RefCOCO(+g) [27, 45], PhraseCut [42], and our proposed MRES-32M. Part-Level: expressions that specifies various parts of the target object in the given image.

Image Source	ReferIt	RefCOCO(+g)	PhraseCut	MRES-32M
	CLEF [11]	COCO [24]	VG [21]	Object365 [35]
Object-Level	✓	✓	✓	✓
Part-Level	×	×	×	✓
Expression Type	free	free	templated	free

[13], various multimodal frameworks such as [13, 19, 41, 43] towards precise RES have been proposed to deal with the challenging feature extraction and alignment problems between visual and linguistic modalities. However, current works are limited to the scope of object-level grounding. As shown in Table 1, the RefCOCO dataset [45] stands as one of the most widely used grounding benchmarks by far, basically containing only object-level visual and textual annotations. *It does not take into account the part-level grounding task, which is crucial for future multimodal models to act as intelligent agents to realize the fine-grained perception of real world.* The existing grounding datasets couldn’t support the training and evaluation of such advanced capabilities. Thus, exploring how to transcend the current object-level constraints and delve into finer-grained part-level grounding, is highly meaningful and worthy of in-depth study, which is the primary focus of this work.

In fact, prior to this work, a few works along another research line have made significant strides towards more fine-grained visual understanding. The demand for finer-grained visual understanding of objects arouses research community into constructing high-quality part-level dataset. Specifically, recent years have witnessed the introduction of various datasets that provide fine-grained part-level masks and bounding boxes annotations for either objects of specific categories [10, 34, 40] or general categories [4, 12, 26, 33]. However, these datasets essentially correspond to unimodal (*i.e.*, visual) downstream tasks, lacking a deep connection between fine-grained part-level masks and rich textual descriptions. To our knowledge, few of previous studies have established this connection. Consequently, there appears to be a limited availability of fine-grained, large-scale vision-language datasets that facilitate part-level cross-modality understanding, which is necessary in terms of two aspects. Firstly, when describing an object, people always naturally gravitate towards detailing part-level local features, underscoring the indispensable need for multimodal agents to understand the part granularity. Secondly, a fine-grained understanding at the part level can positively promote the perception of object-level targets, especially under extreme conditions such as occlusion or deformation, which will consistently propel advancements in widely focused object-level tasks like classic visual grounding.

Therefore, in this work, we attempt to fill this impor-

tant blank space that has been neglected before and move towards finer-grained part-level RES task. Specifically, to push towards finer-grained vision-language understanding, we propose a new multi-granularity referring expression segmentation (MRES) task and an evaluation benchmark named RefCOCO_M by manually annotating the part-level targets based on the previous benchmark RefCOCO with only object-level labels. We also construct the largest-scale visual grounding dataset which is also the first dataset to support part-level visual-textual annotations, and build a simple baseline model to accomplish the unified multi-granularity (*i.e.*, object-level and part-level) RES task.

Our main contributions can be summarized as follows:

- We propose a new MRES task (as shown in Fig. 1) with corresponding benchmark RefCOCO_M for evaluation, pushing classic RES task towards finer-granularity understanding of real-world scenes.
- We build a multi-grained visual grounding dataset namely MRES-32M, which to the best of our knowledge is the first grounding dataset that supports part-level vision-language annotations and also the largest-scale visual grounding dataset.
- To effectively unify both the object-level and part-level RES tasks, we propose a simple yet strong model namely UniRES, which achieves the new state-of-the-art (SOTA) performance on three object-level benchmarks for classic RES tasks and our multi-granularity RefCOCO_M benchmark for the proposed MRES task.

2. Related Work

Referring Expression Segmentation. As a challenging visual grounding task the concept of RES is first proposed by [13]. Subsequent works such as [15, 16, 41, 44, 46] predominantly follow a two-step pipeline of encoding the linguistic and visual features separately and deriving the fused multimodal features from unimodal representations for mask prediction. In this way, the effectiveness of obtained multimodal representations essentially dominates the model performance, which has been continually studied in the following research [7, 22, 43, 49, 50]. Recently, there are some works [28, 37, 47] focusing on zero-shot RES task which has great application potential. In addition, a few recent works [14, 25] have been devoted to address the limitations of the existing benchmark datasets [18, 27, 42, 45] for visual grounding task. However, it’s worth noting that previous works primarily focus on classic object-level RES methods and datasets, paying few attention to the vital fine-grained part-level grounding.

Part-level Visual Perception. The increasing interest in the fine-grained understanding of objects has driven the creation of part-level annotated datasets across both specialized and general domains. Pioneering research in the former field has introduced datasets with part-level annotations, fo-

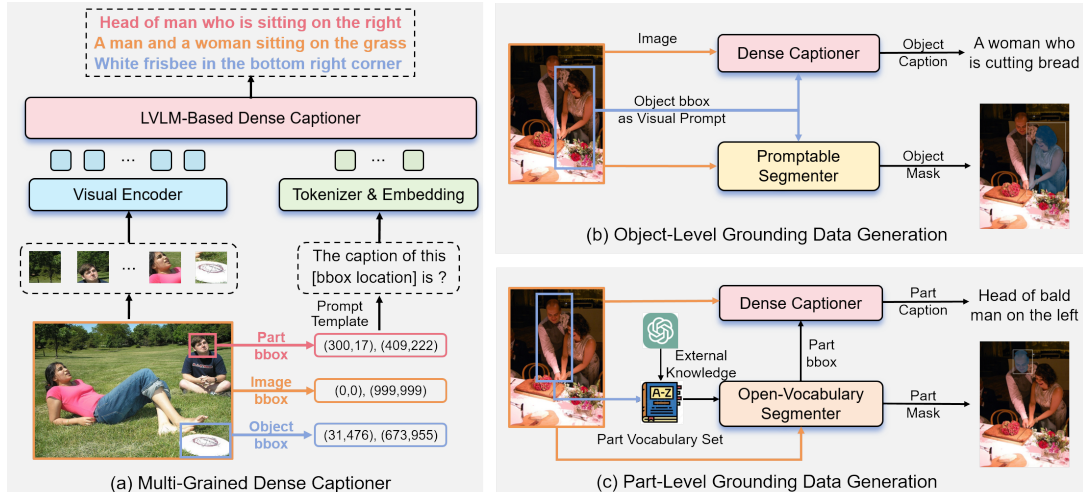


Figure 3. The illustration of our data engine for building the MRES-32M dataset. (a) We start by fine-tuning an LVLm to create a capable dense captioner, which can effectively handle captioning tasks at three levels of granularity. (b) To generate object-level grounding data, we feed images and original bounding boxes into the dense captioner and a powerful segmenter to obtain the captions and masks of various objects. (c) We leverage the external knowledge from LLMs to decompose the existing object category annotations into a vocabulary set of part-level tags, which are sequentially fed into an open-vocabulary segmenter and our captioner to acquire the part-level annotations.

The constrained training data have resulted in the fine-grained descriptive capabilities of LVLms being confined merely to the object-level. To fully harness the open-world visual knowledge acquired from extensive pretraining on image-level and object-level data, we have devised a unified fine-tuning scheme that equips LVLms to operate as dense captioners across all granularity levels, as depicted in Fig. 3 (a). Data for all three granularity levels are sourced from manual annotations to ensure reliability. The input comprises an image and the corresponding bounding box to be described, with bounding box coordinates normalized to integer values within the $[0,999]$ range. For image granularity, we employ the COCO dataset [24], where all bboxes are uniformly represented as $(0,0), (999,999)$. For object granularity, we utilize the Visual Genome dataset [21]. For part granularity, we draw upon unimodal semantic segmentation data [4, 12, 33] and employ a template in the form of *PartNameX of ObjectNameY* to construct dense captions. This unified multitask training approach can be synergistic across different granularity: it allows LVLms to incorporate more comprehensive and detailed information to enhance part granularity descriptions. Concurrently, knowledge of part granularity assists LVLms in generalizing knowledge within object interiors.

Model-Assisted Data Generation. For the data generation of object-level visual grounding, we capitalize on large-scale object detection dataset Object365 [35] to furnish highly reliable bounding boxes. Moreover, the rich category labels inherent in it ensure a comprehensive knowledge coverage. As illustrated in Fig. 3 (b), the bounding boxes will serve as visual prompts, which are independently sent into both a promptable segmenter (*i.e.* segment anything model

[20]) and our dense captioner to obtain the segmentation masks and detailed semantic descriptions, respectively.

For part-level grounding data, we propose a hierarchical annotation scheme that builds upon existing object-level annotations, as shown in Fig. 3 (c). Specifically, employing GPT-4 [29] endowed with extensive external knowledge, we decompose the objects present in a given image to generate a customized part vocabulary set. This tailored vocabulary set is then fed into an open-vocabulary segmenter [36], which yields precise part masks and bounding boxes. These bounding boxes are subsequently sent into our dense captioner to acquire corresponding detailed captions.

Data Filtering. After completing the multi-granularity annotation of all images, we further introduce CLIP [31] for filtering. The bounding box is cropped from the original image, and then sent to the encoder together with the dense caption to measure the similarity. To ensure the consistency between the visual and linguistic annotations to a great extent, we retain box-caption pairs with similarity greater than 0.5 as the final annotation results.

4.2. MRES-32M Dataset Details

As listed in Table 2, existing datasets, such as the most commonly used benchmark dataset RefCOCO [45], have limitations in terms of small data scale and lacking part-level dense annotations. As one of the pioneering works for finer-granularity unimodal visual comprehension, PACO [33] only contains semantic tags of objects or parts, without informative descriptions and visual context encompassed as referring expressions. Summarily, we compare the proposed MRES-32M with existing datasets and list some unique and significant properties of our dataset in Table

Table 2. Comparisons with previous object-level visual grounding datasets and part-level segmentation datasets. # denotes the specific number, where Cats and Avg Len denote the object/part categories and the average length of referring expressions. “-” denotes the corresponding part masks or captions are unavailable.

Dataset	#Imgs	#Objs	#Parts	#Cats	#Avg Len
Object-Level Visual Grounding					
ReferIt [18]	20K	97K	-	238/-	3.2
RefCOCO [45]	20K	50K	-	80/-	3.6
RefCOCO+ [45]	20K	49K	-	80/-	3.5
RefCOCOg [27]	26K	54K	-	80/-	8.4
GRES [25]	20K	60K	-	80/-	3.7
Part-Level Segmentation & Detection					
PartsIN [12]	24K	24K	112K	158/609	-
PascalPart [4]	19K	40K	363K	20/193	-
PACO [33]	20K	260K	641K	75/456	-
Multi-Grained Visual Grounding					
MRES-32M (Ours)	1M	15.3M	16.9M	365/2299	4.6

2. Besides, we have also provided a few examples in our MRES-32M dataset in the appendix.

Unified Multi-Granularity. In comparison to the grounding counterparts, our MRES-32M is the first visual grounding dataset covering both part and object granularity. In comparison to the part-level segmentation counterparts, our MRES-32M provides informative and unique fine-grained description for each part mask.

More Diversified Categories. Our MRES-32M is composed of 365 object categories and an associated 2,299 part categories. Compared with existing datasets, it covers a wider range of multimodal knowledge and is an important step towards open-world understanding.

Breakable Data Scales. To the best of our knowledge, MRES-32M is the largest dataset in the current grounding research community. In terms of the number of images and object instances, it surpasses the largest existing visual grounding dataset RefCOCOg [27] by factors of 38 and 283, respectively. Meanwhile, it encompasses part instance counts that exceed the largest existing part semantic segmentation dataset [33] by 58 times.

More Complex References. Benefiting from our LVLM-based dense captioner, the reference of MRES-32M could be more fully combined with the visual context for entity (*i.e.*, part and object) description. Without sticking to a specific template, the relationships and attributes of entities could be highlighted in free natural language expressions.

5. Multi-Granularity RES Model

Next, we describe the proposed multi-granularity RES model UniRES for the unified MRES task with the referring targets at both object and part granularity. *Since our original intention is to establish a simple and easy-to-follow baseline model for the proposed multi-granularity RES task,*

the structure of our model UniRES is designed to be simple and clean. As shown in Fig. 4, UniRES has three major components, which will be illustrated below.

Vision-Language Backbone. Taking account of the required attributes of both strong ability to capture the vision-language feature representations and promising scalability, we leverage the CLIP pre-trained weights of CLIP model [32], which learns transferable visual and linguistic concepts from vast image-text pairs, and adopt it as the backbone for our referring segmentation framework. The respective image and text encoders (*i.e.* Vision Transformer (ViT) [8] and Transformer [39] respectively) from CLIP are utilized to effectively extract visual and linguistic features.

Query-based Grouping Design. To effectively accomplish the proposed MRES task, it is essential to exploit both low-level local and high-level global visual features. In order to enhance the local-global visual representations of CLIP backbone without introducing much additional parameters or altering model’s structure, we have incorporated 64 and 8 learnable tokens (set empirically) into the first and middle layers of CLIP visual backbone. These learnable tokens traverse the first and second halves of the visual backbone. We expect that the ViT’s internal self-attention mechanism implicitly serves as a manner to perform visual grouping, obtaining representative group tokens that capture the low-level local and high-level global features simultaneously. Based on the fact that local features spread more fragmented, the number of appended low-level group tokens is greater than that of high-level counterparts. Then, group tokens from both levels are fed into a language-guided region filter (LRF) to select the language-related visual features by cross-attention mechanism, which is followed by concatenation to fuse these expression-related visual group tokens for subsequent vision-language decoding.

Two-Stage V-L Decoder. Now that the visual and textual feature representations from backbone, and the expression-related group tokens with two different levels are obtained, the two-stage mask decoder which comprises stacked Transformer layers is utilized to generate the segmentation masks. Specifically, the first-stage V-L decoder takes the extracted visual and textual features as input and generates first-stage fused multimodal representations. Subsequently, these multimodal features are further integrated with the grouped expression-related region features of the two semantic levels (*i.e.*, low-level & high-level) to realize further feature enhancement, which is followed by a linear projection layer to obtain the final segmentation masks.

6. Experimental Results

To evaluate the effectiveness of our method, comprehensive experiments are conducted on three classic RES datasets (*i.e.*, RefCOCO [45], RefCOCO+ [45], RefCOCOg [27]) and our RefCOCO_m benchmark for multi-granularity RES.

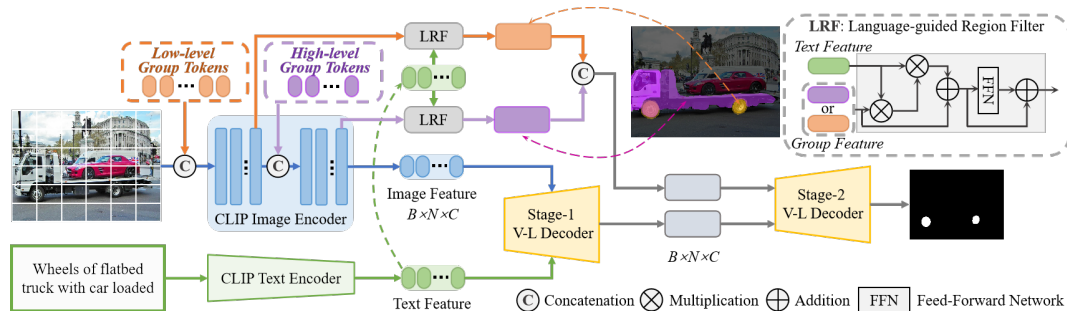


Figure 4. The architecture of our UniRES model as a simple baseline for the MRES task. UniRES mainly comprises three parts: the visual and textual backbone for feature extraction, the pixel grouping design for aggregating the low-level and high-level features, and the cascaded two-stage vision-language decoder for multimodal feature fusion and the generation of segmentation masks.

6.1. Datasets

The details about the three benchmark datasets for classic RES task can be found in the appendix.

6.2. Experimental Setup

Implementation Details. The implementation details are provided in the appendix.

Evaluation metrics. We adopt mean Intersection-over-Union (mIoU) and overall Intersection-over-Union (oIoU) as evaluation metrics. The mIoU measures the ratio between the intersection area and the union area of the predictions and labels among the test samples, while the oIoU calculates the total intersection area over total union area.

6.3. Main Results

6.3.1 Multi-Granularity MRES Task

Comparison with SOTA Methods. To evaluate the multi-granularity grounding performance of our UniRES as the baseline for the proposed MRES task and previous RES methods, we further conduct experimental comparison on our newly built RefCOCO_m benchmark dataset. As presented in Table 3, the classic RES methods including four specialist models (*i.e.*, SeqTR [48], CRIS [41], LAVT [43]) and two generalist models (*e.g.*, X-Decoder [49] and SEEM [50]) are incorporated. For fair comparisons, we re-implement these SOTA methods and report their performance on our RefCOCO_m. It is clear that either the specialist models for classic RES task or the powerful generalist models perform poorly on RefCOCO_m, which requires the crucial skills of both part-level and object-level referring segmentation. Benefiting from our MRES-32M dataset, our UniRES can better master the part-level RES skills and handle the multi-granularity RES task, achieving considerably higher segmentation accuracy. Due to the significantly higher difficulty of part-level RES (*i.e.* multi-granularity RES) compared to classic RES, the absolute value of segmentation accuracy is correspondingly lower. This further emphasizes the importance of researching finer-grained part grounding where previous SOTA methods have fallen short.

Table 3. Comparison with previous SOTA methods on our RefCOCO_m benchmark in terms of mIoU. Part and Obj&Part denote part-only and multi-grained evaluation settings of our MRES task.

Methods	val		testA		testB	
	Part	Obj&Part	Part	Obj&Part	Part	Obj&Part
<i>Specialists</i>						
SeqTR [48]	13.9	28.2	12.1	22.8	18.1	34.7
CRIS [41]	10.6	25.4	10.1	21.2	12.9	30.0
LAVT [43]	15.3	29.9	13.2	24.4	18.7	35.5
<i>Generalists</i>						
X-Decoder [49]	16.2	29.5	13.6	23.6	20.3	33.8
SEEM [50]	16.1	29.4	13.6	23.4	20.4	33.9
UniRES (Ours)	19.6	34.3	16.4	27.8	25.2	41.7

Qualitative Analysis. We also conduct visual comparisons of CRIS [41], LAVT [43] and our UniRES on the MRES task, which is provided in appendix.

6.3.2 Classic Object-Level RES Task

Comparison with SOTA Methods. To validate the superiority of our MRES-32M dataset and model UniRES, our framework is fairly evaluated against previous SOTA methods on RefCOCO [45], RefCOCO+ [45] and RefCOCO_g [27] under both zero-shot and fine-tuning settings. As presented in Table 4, our method greatly outperforms previous methods in terms of segmentation accuracy across all the benchmark datasets. Via directly zero-shot transferring to the downstream classic RES task after pre-training on our MRES-32M dataset, Our UniRES model achieves a leading zero-shot segmentation accuracy (*i.e.*, approximately 71%) without any fine-tuning to adapt to downstream task’s data, which is significantly higher than all the recently proposed zero-shot RES methods by a large margin (*i.e.*, $\uparrow 30\text{-}40\%$ mIoU). At the same time, it is worth noting that our zero-shot segmentation performance is already better than many of previous RES methods under fine-tuning setting (*e.g.*, CRIS [41] and ReSTR [19]), validating the potential of our MRES-32M dataset and UniRES model. Furthermore, after fine-tuning on the classic RES datasets, our UniRES greatly surpasses previous methods no matter the specialist models (*e.g.*, LAVT [43] and LISA [22]) for classic RES or the generalist models (*e.g.*, X-Decoder [49] and SEEM [50]).

Table 4. Comparisons with the state-of-the-art approaches on previous three classic RES benchmark datasets under both the zero-shot and fine-tuning settings. “-” denotes that the result is not provided.

Method		RefCOCO			RefCOCO+			RefCOCOg	
		val	testA	testB	val	testA	testB	val	test
<i>Zero-Shot Methods</i>									
mIoU	Region token [47] (CVPR-23)	23.4	22.1	24.6	24.5	22.6	25.4	27.6	27.3
	Cropping [47] (CVPR-23)	24.8	22.6	25.7	26.3	24.1	26.5	31.9	30.9
	Global-Local CLIP [47] (CVPR-23)	26.2	24.9	26.6	27.8	25.6	27.8	33.5	33.7
	SAM-CLIP [28] (arXiv-23)	26.3	25.8	26.4	25.7	28.0	26.8	38.8	38.9
	Ref-Diff [28] (arXiv-23)	37.2	38.4	37.2	37.3	40.5	33.0	44.0	44.5
	TAS [37] (arXiv-23)	39.8	41.1	36.2	43.6	49.1	36.5	46.6	46.8
	UniRES (Ours)	71.2	74.8	66.0	59.9	66.7	51.4	62.3	63.2
<i>Fine-Tune Methods</i>									
oIoU	EFNet [9] (CVPR-21)	62.8	65.7	59.7	51.5	55.2	43.0	-	-
	LTS [17] (CVPR-21)	65.4	67.8	63.1	54.2	58.3	48.0	54.4	54.3
	ReSTR [19] (CVPR-22)	67.2	69.3	64.5	55.8	60.4	48.3	-	-
	ReLA [25] (CVPR-23)	73.8	76.5	70.2	66.0	71.0	57.7	65.0	66.0
	X-Decoder [49] (CVPR-23)	-	-	-	-	-	-	64.6	-
	SEEM [50] (arXiv-23)	-	-	-	-	-	-	65.7	-
	LISA [22] (arXiv-23)	74.9	79.1	72.3	65.1	70.8	58.1	67.9	70.6
	UniRES (Ours)	77.4	80.9	74.7	69.4	76.1	61.4	69.0	71.7
	mIoU	VLT [7] (ICCV-21)	65.7	68.3	62.7	55.5	59.2	49.4	53.0
RefTr [23] (NeurIPS-21)		74.3	76.8	70.9	66.8	70.6	59.4	66.6	67.4
SeqTR [48] (ECCV-22)		71.7	73.3	69.8	63.0	66.7	59.0	65.0	65.7
CRIS [41] (CVPR-22)		70.5	73.2	66.1	62.3	68.1	53.7	59.9	60.4
LAVT [43] (CVPR-22)		74.5	76.9	70.9	65.8	71.0	59.2	63.3	63.6
UniRES (Ours)		79.2	81.6	76.6	73.0	78.1	65.8	71.7	73.2

6.4. Ablation Studies

We conduct ablation experiments on our RefCOCO_m validation set. *The tables below involves three evaluation settings on our RefCOCO_m, using only objects as RES targets, using only parts as RES targets, and a mixed setting that combines both of the above two granularity levels.*

6.4.1 Ablation Study on MRES-32M Dataset

Data Granularity. We first probe into the effect of pre-training data’s granularity with 50% of our MRES-32M dataset. As shown in Table 5, the baseline without pre-training on MRES-32M dataset obtains 75.2%, 15.8% and 32.0% mIoU separately on our RefCOCO_m validation set under three granularity levels. Either introducing the object-level or part-level data from MRES-32M dataset for pre-training consistently results in an considerable accuracy increase across different granularity settings. In fact, pre-training on data at a certain granularity level to improve performance on corresponding benchmark aligns with common sense. However, it is noteworthy that incorporating part-level data into training also enhances model performance on the only object-level RefCOCO_m. This underscores the significance of part-level understanding, as emphasized in the Sec. 1, where grasping the nuances of parts can yield benefits for object-level grounding. Besides, by

jointly incorporating the training samples of both granularity, our method attains 2.8% improvement against baseline under all the granularity settings on RefCOCO_m, which fully demonstrate the benefit of exploiting our MRES-32M dataset for both object-level and part-level RES tasks.

Table 5. Ablation study on the granularity of pre-training data in our proposed MRES-32M. Object and Part denote the introduction of object-level and part-level data for pre-training.

Object	Part	RefCOCO _m		
		Object-Only	Part-Only	Object & Part
		75.2	15.8	30.5
✓		77.5	15.9	31.1
	✓	75.9	18.4	32.6
✓	✓	78.0	18.6	33.3

Table 6. Ablation study on the data scale of MRES-32M dataset.

Ratios	RefCOCO _m		
	Object-Only	Part-Only	Object & Part
0%	75.2	15.8	30.5
20%	76.8	17.3	32.0
50%	78.0	18.6	33.3
100%	79.2	19.6	34.3

Data Scale. Next, we investigate the effect of different percentages of training samples in our MRES-32M dataset. The results are presented in Table 6. It is obvious in Table 6 that the model performance for our MRES tasks is

consistently improved with more and more employed training samples, which implicitly validates the high quality of our built dataset. As the ratios of introduced training samples continue to rise, there’s no sign of diminishing gains in model performance, suggesting that our framework has significant potential with continually scaled up training data.

Table 7. Ablation study on the effectiveness of our MRES-32M dataset on RES SOTA methods.

Methods	MRES-32M	RefCOCO _m		
		Object-Only	Part-Only	Object & Part
CRIS		70.5	10.6	25.4
CRIS	✓	73.1	15.5	29.7
LAVT		74.5	15.3	29.9
LAVT	✓	75.7	19.3	33.2

Dataset Necessity. To prove the necessity and effectiveness of our MRES-32M dataset for the proposed MRES task, we take previous SOTA methods (*i.e.*, CRIS [41] and LAVT [43]) on classic RES and compare the segmentation accuracy of the same models with or without pre-training on our MRES-32M. As shown in Table 7, the original CRIS and LAVT can already well handle the classic RES task with only object-level grounding skills, but it performs poorly on the multi-grained RefCOCO_m benchmark, which requires part-level grounding capability. In contrast, pre-training on our MRES-32M consistently leads to great performance gains no matter on object-only, part-only or the multi-granularity RES tasks, because our high-quality MRES-32M dataset can effectively enable the model to handle the part-level visual grounding task and enhances the original object-level grounding capability.

6.4.2 Ablation Study on UniRES Model

Table 8. Ablation study on the query-based grouping design.

High-Level	Low-Level	RefCOCO _m		
		Object-Only	Part-Only	Object & Part
		74.4	14.9	29.6
✓		74.9	15.2	30.0
	✓	74.9	15.4	30.1
✓	✓	75.2	15.8	30.5

Besides, we additionally conduct ablation studies on the structural design of our model UniRES. As presented in Table 8, removing the simple but effective grouping design at two different levels (this simultaneously leads to discarding the second-stage decoder) will result in a considerable deterioration in model performance. Since the sequentially appended pixel grouping tokens are introduced to effectively capture the high-level and low-level visual features for further visual feature enhancement, losing any type of these group tokens will lead to a decrease in the model’s segmentation accuracy on both object-level and part-level RES tasks. In order to verify that the high-level and low-level pixel grouping tokens can capture the clustering features of the corresponding level, we also visualize the group tokens

of the two levels (before being sent to LRF) for qualitative analysis, where the different colors represent the various areas of the same clustering group. The visualization results in Fig. 5 affirms that our introduced low-level and high-level group tokens respectively capture the part-level local features and aggregate the object-level global features with stronger semantics through the pixel grouping process, which aligns with our intentions for the grouping design.

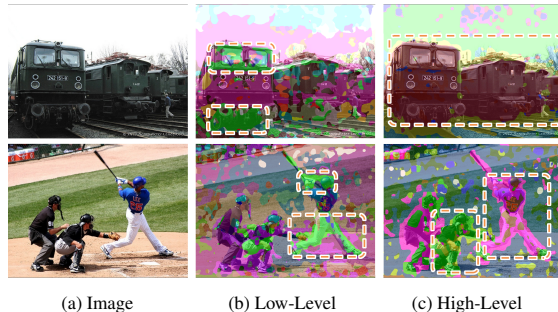


Figure 5. Qualitative analysis for ablation study on the object-level and part-level grouping design in our model structure. (a) the input image. (b) low-level group tokens. (c) high-level group tokens.

7. Conclusion and Broader Impact

In this paper, we move beyond previous works that focused solely on object-level visual grounding tasks and take a step further to finer-grained part-level RES. We put forward a new multi-granularity referring expression segmentation task and establish an evaluation benchmark named RefCOCO_m by manual annotation. To advance the visual grounding at both object and part levels towards finer-grained vision-language understanding, we build the largest visual grounding dataset MG-32M to date, which is also the first dataset to provide part-level vision-language annotations. Furthermore, we have developed a simple yet strong multi-grained referring segmentation model called UniRES. As a baseline for our newly proposed MRES task, UniRES achieves the new state-of-the-art performance on both our RefCOCO_m for MRES task and three classic RES datasets. We plan to release our RefCOCO_m benchmark, the MG-32M dataset, and the UniRES model to the public, aspiring to foster future research in fine-grained visual grounding tasks and to inspire new research in this direction.

Acknowledgement

We thank Yepeng Tang for the helpful discussions on this work, all the Image & Video Analysis Group (IVA)’s members in CASIA for the technical support, and all the insightful reviewers for the helpful suggestions. This work was supported by the National Science and Technology Major Project (No.2022ZD0118801), National Natural Science Foundation of China (U21B2043, 62206279).

References

- [1] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023. [3](#)
- [2] Jun Chen, Deyao Zhu¹ Xiaoqian Shen¹ Xiang Li, Zechun Liu² Pengchuan Zhang, Raghuraman Krishnamoorthi² Vikas Chandra² Yunyang Xiong, and Mohamed Elhoseiny. Minigt-v2: Large language model as a unified interface for vision-language multi-task learning. *arXiv preprint arXiv:2310.09478*, 2023.
- [3] Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. Shikra: Unleashing multi-modal llm’s referential dialogue magic. *arXiv preprint arXiv:2306.15195*, 2023. [3](#)
- [4] Xianjie Chen, Roozbeh Mottaghi, Xiaobai Liu, Sanja Fidler, Raquel Urtasun, and Alan Yuille. Detect what you can: Detecting and representing objects using holistic models and body parts. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1971–1978, 2014. [2](#), [3](#), [4](#), [5](#)
- [5] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. See <https://vicuna.lmsys.org> (accessed 14 April 2023), 2023. [3](#)
- [6] Daan de Geus, Panagiotis Meletis, Chenyang Lu, Xiaoxiao Wen, and Gijs Dubbelman. Part-aware panoptic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5485–5494, 2021. [3](#)
- [7] Henghui Ding, Chang Liu, Suchen Wang, and Xudong Jiang. Vision-language transformer and query generation for referring segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16321–16330, 2021. [2](#), [7](#)
- [8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. [5](#), [1](#)
- [9] Guang Feng, Zhiwei Hu, Lihe Zhang, and Huchuan Lu. Encoder fusion network with co-attention embedding for referring image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15506–15515, 2021. [7](#)
- [10] Ke Gong, Xiaodan Liang, Dongyu Zhang, Xiaohui Shen, and Liang Lin. Look into person: Self-supervised structure-sensitive learning and a new benchmark for human parsing. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 932–940, 2017. [2](#), [3](#)
- [11] Michael Grubinger, Paul Clough, Henning Müller, and Thomas Deselaers. The iapr tc-12 benchmark: A new evaluation resource for visual information systems. In *International workshop ontoImage*, 2006. [2](#)
- [12] Ju He, Shuo Yang, Shaokang Yang, Adam Kortylewski, Xiaoding Yuan, Jie-Neng Chen, Shuai Liu, Cheng Yang, Qihang Yu, and Alan Yuille. Partimagenet: A large, high-quality dataset of parts. In *European Conference on Computer Vision*, pages 128–145. Springer, 2022. [2](#), [3](#), [4](#), [5](#)
- [13] Ronghang Hu, Marcus Rohrbach, and Trevor Darrell. Segmentation from natural language expressions. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part I 14*, pages 108–124. Springer, 2016. [2](#), [1](#)
- [14] Yutao Hu, Qixiong Wang, Wenqi Shao, Enze Xie, Zhenguo Li, Jungong Han, and Ping Luo. Beyond one-to-one: Rethinking the referring image segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4067–4077, 2023. [2](#)
- [15] Zhiwei Hu, Guang Feng, Jiayu Sun, Lihe Zhang, and Huchuan Lu. Bi-directional relationship inferring network for referring image segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4424–4433, 2020. [2](#)
- [16] Shaofei Huang, Tianrui Hui, Si Liu, Guanbin Li, Yunchao Wei, Jizhong Han, Luoqi Liu, and Bo Li. Referring image segmentation via cross-modal progressive comprehension. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10488–10497, 2020. [2](#)
- [17] Ya Jing, Tao Kong, Wei Wang, Liang Wang, Lei Li, and Tieniu Tan. Locate then segment: A strong pipeline for referring image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9858–9867, 2021. [7](#)
- [18] Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. Referitgame: Referring to objects in photographs of natural scenes. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 787–798, 2014. [2](#), [5](#)
- [19] Namyup Kim, Dongwon Kim, Cuiling Lan, Wenjun Zeng, and Suha Kwak. Restr: Convolution-free referring image segmentation using transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18145–18154, 2022. [2](#), [6](#), [7](#)
- [20] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023. [3](#), [4](#)
- [21] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*, 123:32–73, 2017. [2](#), [4](#)
- [22] Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. Lisa: Reasoning segmentation via large language model. *arXiv preprint arXiv:2308.00692*, 2023. [2](#), [6](#), [7](#)
- [23] Muchen Li and Leonid Sigal. Referring transformer: A one-step approach to multi-task visual grounding. *Advances*

- in neural information processing systems, 34:19652–19664, 2021. 7
- [24] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014. 2, 4, 1
- [25] Chang Liu, Henghui Ding, and Xudong Jiang. Gres: Generalized referring expression segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23592–23601, 2023. 2, 5, 7
- [26] Kaichun Mo, Shilin Zhu, Angel X Chang, Li Yi, Subarna Tripathi, Leonidas J Guibas, and Hao Su. Partnet: A large-scale benchmark for fine-grained and hierarchical part-level 3d object understanding. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 909–918, 2019. 2, 3
- [27] Varun K Nagaraja, Vlad I Morariu, and Larry S Davis. Modeling context between objects for referring expression understanding. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*, pages 792–807. Springer, 2016. 2, 5, 6, 1
- [28] Minheng Ni, Yabo Zhang, Kailai Feng, Xiaoming Li, Yiwen Guo, and Wangmeng Zuo. Ref-diff: Zero-shot referring image segmentation with generative models. *arXiv preprint arXiv:2308.16777*, 2023. 2, 7
- [29] R OpenAI. Gpt-4 technical report. arxiv 2303.08774. *View in Article*, 2023. 4
- [30] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019. 1
- [31] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 4
- [32] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 5, 1
- [33] Vignesh Ramanathan, Anmol Kalia, Vladan Petrovic, Yi Wen, Baixue Zheng, Baishan Guo, Rui Wang, Aaron Marquez, Rama Kovvuri, Abhishek Kadian, et al. Paco: Parts and attributes of common objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7141–7151, 2023. 2, 3, 4, 5
- [34] N Dinesh Reddy, Minh Vo, and Srinivasa G Narasimhan. Carfusion: Combining point tracking and part detection for dynamic 3d reconstruction of vehicles. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1906–1915, 2018. 2, 3
- [35] Shuai Shao, Zeming Li, Tianyuan Zhang, Chao Peng, Gang Yu, Xiangyu Zhang, Jing Li, and Jian Sun. Objects365: A large-scale, high-quality dataset for object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 8430–8439, 2019. 2, 4
- [36] Peize Sun, Shoufa Chen, Chenchen Zhu, Fanyi Xiao, Ping Luo, Saining Xie, and Zhicheng Yan. Going denser with open-vocabulary part segmentation. *arXiv preprint arXiv:2305.11173*, 2023. 3, 4
- [37] Yucheng Suo, Linchao Zhu, and Yi Yang. Text augmented spatial-aware zero-shot referring image segmentation. *arXiv preprint arXiv:2310.18049*, 2023. 2, 7
- [38] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023. 3
- [39] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 5
- [40] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. 2011. 2, 3
- [41] Zhaoqing Wang, Yu Lu, Qiang Li, Xunqiang Tao, Yandong Guo, Mingming Gong, and Tongliang Liu. Cris: Clip-driven referring image segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11686–11695, 2022. 2, 6, 7, 8, 1
- [42] Chenyun Wu, Zhe Lin, Scott Cohen, Trung Bui, and Subhransu Maji. Phrasecut: Language-based image segmentation in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10216–10225, 2020. 2
- [43] Zhao Yang, Jiaqi Wang, Yansong Tang, Kai Chen, Hengshuang Zhao, and Philip HS Torr. Lavt: Language-aware vision transformer for referring image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18155–18165, 2022. 2, 6, 7, 8, 1
- [44] Linwei Ye, Mrigank Rochan, Zhi Liu, and Yang Wang. Cross-modal self-attention network for referring image segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10502–10511, 2019. 2
- [45] Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. Modeling context in referring expressions. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part II 14*, pages 69–85. Springer, 2016. 2, 3, 4, 5, 6, 1
- [46] Licheng Yu, Zhe Lin, Xiaohui Shen, Jimei Yang, Xin Lu, Mohit Bansal, and Tamara L Berg. Mattnet: Modular attention network for referring expression comprehension. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1307–1315, 2018. 2

- [47] Seonghoon Yu, Paul Hongsuck Seo, and Jeany Son. Zero-shot referring image segmentation with global-local context features. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19456–19465, 2023. [2](#), [7](#)
- [48] Chaoyang Zhu, Yiyi Zhou, Yunhang Shen, Gen Luo, Xingjia Pan, Mingbao Lin, Chao Chen, Liujuan Cao, Xiaoshuai Sun, and Rongrong Ji. Seqtr: A simple yet universal network for visual grounding. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXV*, pages 598–615. Springer, 2022. [6](#), [7](#)
- [49] Xueyan Zou, Zi-Yi Dou, Jianwei Yang, Zhe Gan, Linjie Li, Chunyuan Li, Xiyang Dai, Harkirat Behl, Jianfeng Wang, Lu Yuan, et al. Generalized decoding for pixel, image, and language. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15116–15127, 2023. [2](#), [6](#), [7](#)
- [50] Xueyan Zou, Jianwei Yang, Hao Zhang, Feng Li, Linjie Li, Jianfeng Gao, and Yong Jae Lee. Segment everything everywhere all at once. *arXiv preprint arXiv:2304.06718*, 2023. [2](#), [6](#), [7](#)

A. Appendix

In this appendix section, we provide following items:

- (Sec. 1) Potential limitations and future works.
- (Sec. 2) Detailed information about the three benchmark datasets (*i.e.* RefCOCO, RefCOCO+ and RefCOCO+) for classic referring expression segmentation (RES) task.
- (Sec. 3) Implementation details on the three benchmark datasets (*i.e.* RefCOCO, RefCOCO+ and RefCOCO+) for classic RES task and the newly built RefCOCO_m benchmark for our multi-granularity RES (MRES) task.
- (Sec. 4) More quantitative analysis about our newly built MRES-32M dataset.
- (Sec. 5) More visualization results of visual comparisons, and the selected samples from our newly built RefCOCO_m benchmark and MRES-32M dataset.

A.1. Limitation and Future Work

One potential limitation of this work is that the data scale of our MRES-32M and model capacity of UniRES could be scaled up further to push SOTA performance. Moreover, in this work we mainly focus on the (M)RES task with visual masks as output, which means that although our framework demonstrates the new crucial part-level referring segmentation capabilities, it currently can not produce textual responses and thus can not handle the tasks related to vision-language conversations. However, the integration of large language models could enhance our framework’s text comprehension and generation abilities, allowing it to be accordingly adjusted to overcome this limitation. This opens up a future research direction to develop a more general and powerful framework based on multimodal large language models that could interact with user-provided textual and visual inputs across multiple levels of granularity.

A.2. Details about the Benchmark Datasets

RefCOCO [45], stands as one of the largest and frequently utilized datasets derived from MSCOCO [24] for the task of referring expression segmentation. It comprises 142,209 annotated expressions with an average expression length of 3.6 words, labeling 50,000 objects across 19,994 images. The dataset is partitioned into 120,624 training samples, 10,834 validation samples, and two test subsets—test A and test B—containing 5,657 and 5,095 instances, respectively.

RefCOCO+ [45] encompasses 141,564 language expressions with a slightly shorter average expression length of 3.5 words, targeting 49,856 objects within 19,992 images. This dataset follows a similar split as RefCOCO, offering 120,624 training, 10,758 validation, 5,726 test A, and 4,889 test B samples. Unique to RefCOCO+, it omits expressions that use absolute location terms, posing an increased challenge for the classic RES task.

RefCOCO_g [27], serves as the third benchmark dataset and contains 104,560 referring expressions, with an signifi-

cantly longer average length of 8.4 words for 54,822 objects in 26,711 images. The language expressions in this dataset is sourced from Amazon Mechanical Turk, marking a distinction from the previous two datasets. Following previous works, we employ the UMD partition standard [13] for our evaluations in this paper.

A.3. Implementation Details

Experimental Setup. Our work is implemented based on Pytorch [30] and trained with NVIDIA A800 GPUs. Considering the crucial scalability and the ease of implementation, the Vision Transformer [8] is adopted as the image encoder for all the experiments. The text and image encoder are initialized by CLIP [32], while the rest part of model weights are randomly initialized. During training, with 128 and 64 batch size respectively, the AdamW optimizer and a weight decay of $5e-4$ are adopted to pre-train and fine-tune our model for 50 and 15 epochs. With a warm-up strategy for 5-epoch pre-training on our MRES-32M and 1-epoch fine-tuning on the specific downstream grounding dataset, the initial learning rate is set to $1e-5$ with a cosine decay schedule. Following CRIS [41], due to the extra [SOS] and [EOS] tokens, the input sentences are set with a maximum sentence length of 17 for RefCOCO, our RefCOCO_m and RefCOCO+, 22 for RefCOCO_g.

During inference, the predicted results by our method is upsampled back to the original image size and binarized with a threshold of 0.35 to the final segmentation result. Any extra post-processing operations or inference tricks can be exploited to further boost the segmentation accuracy but are not included in this work.

A.4. More Analysis about MRES-32M Dataset

Our MRES-32M is composed of 365 object categories and an associated 2,299 part categories. Compared with existing datasets, it covers a wider range of multimodal knowledge and is an important step towards open-world understanding. Among our MRES-32M dataset, the number of referring expressions per objects’ category and per parts’ category, the word cloud that highlights the head objects’ and parts’ categories are both presented in Fig. 6.

A.5. More Visualization Results

Visual Comparison with SOTA Methods. In addition, to validate the segmentation quality of our framework, classic RES methods CRIS [41] and LAVT [43], as well as our proposed UniRES are further adopted for qualitative comparison. The provided visualization results in the first and second row of Fig. 7 convinces that our method can better accomplish the classic object-level RES task and generate much better fine-grained segmentation masks of the target objects, greatly reducing the over-segmentation and under-segmentation errors. Similarly, as shown in the third and

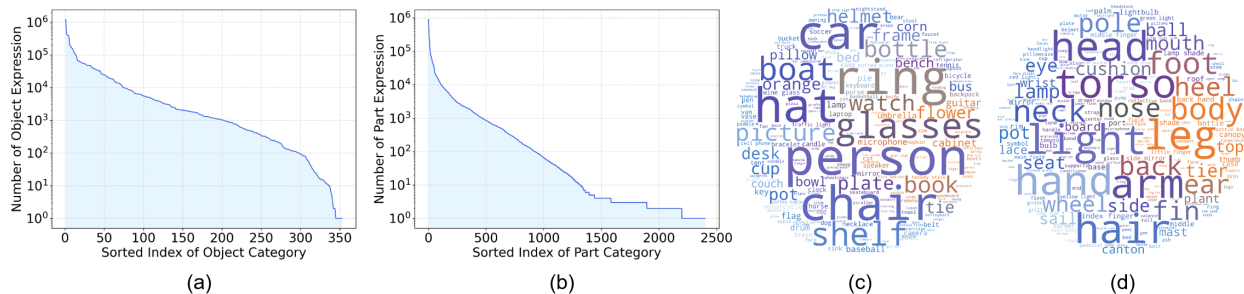


Figure 6. MRES-32M dataset statistics. (a) the number of referring expressions per objects’ category in the log scale. (b) the number of referring expressions per parts’ category in the log scale. (c) the word cloud highlights the head objects’ categories. (d) the word cloud highlights the head parts’ categories.

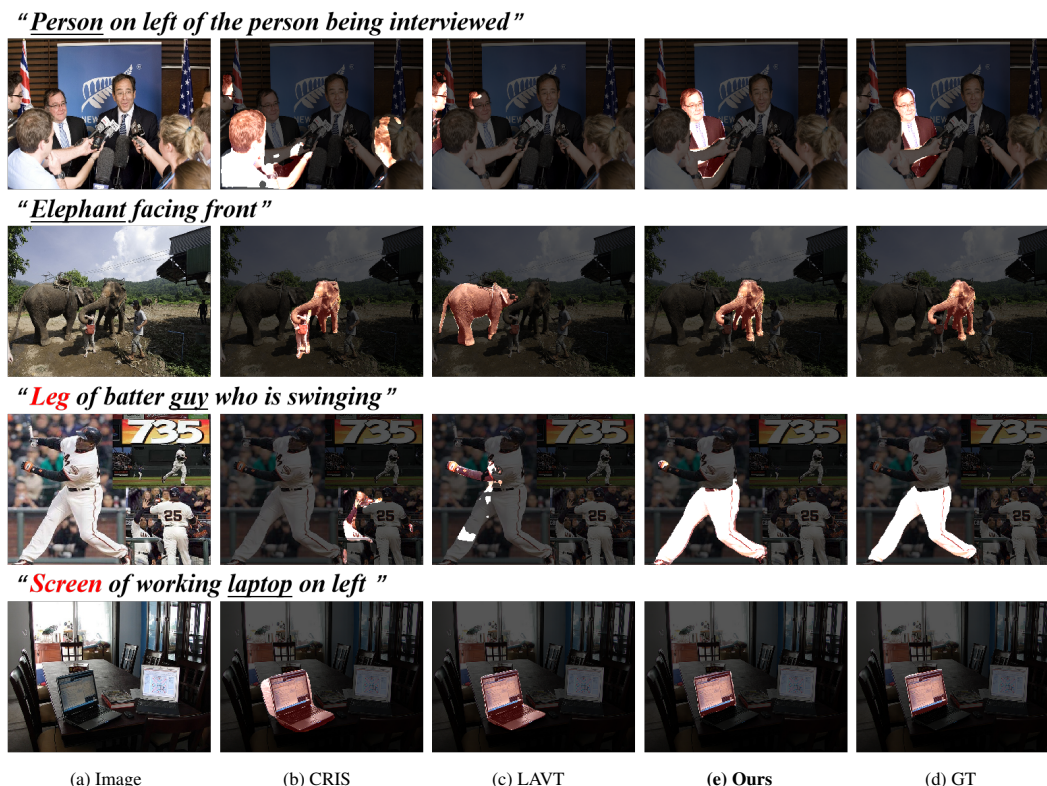


Figure 7. The visual comparison of segmentation results on our RefCOCOM validation set. (a) the input image. (b) CRIS. (c) LAVT. (d) our UniRES. (e) the ground truth.

fourth row in Fig. 7, when facing the more challenging part-level grounding task, our approach can clearly locate and segment the referring target regions more accurately while the other previous methods fail to reach the same level.

Samples of RefCOCOM Benchmark. We have also provided a few more examples in our RefCOCOM benchmark for proposed MRES task in Fig. 8.

Samples of MRES-32M Dataset. A few examples in our newly built MRES-32M dataset for visual grounding task are presented in Fig. 9.



Figure 8. More selected samples from our proposed RefCOCO_m benchmark for multi-granularity RES task.



Figure 9. Selected samples from our built MRES-32M dataset for visual grounding tasks.