

---

# Revisiting Referring Expression Comprehension Evaluation in the Era of Large Multimodal Models

---

**Jierun Chen\***

HKUST

jcheneh@cse.ust.hk

**Fangyun Wei\*†**

Microsoft Research Asia

fawe@microsoft.com

**Jinjing Zhao**

The University of Sydney

jzha0100@uni.sydney.edu.au

**Sizhe Song**

HKUST

ssongad@cse.ust.hk

**Bohuai Wu**

HKUST

bwual@cse.ust.hk

**Zhuoxuan Peng**

HKUST

zpengac@cse.ust.hk

**S.-H. Gary Chan**

HKUST

gchan@cse.ust.hk

**Hongyang Zhang**

University of Waterloo

hongyang.zhang@uwaterloo.ca

## Abstract

1 Referring expression comprehension (REC) involves localizing a target instance  
2 based on a textual description. Recent advancements in REC have been driven by  
3 large multimodal models (LMMs) like CogVLM, which achieved 92.44% accuracy  
4 on RefCOCO. However, this study questions whether existing benchmarks such as  
5 RefCOCO, RefCOCO+, and RefCOCOg, capture LMMs' comprehensive capabilities.  
6 We begin with a manual examination of these benchmarks, revealing high  
7 labeling error rates: 14% in RefCOCO, 24% in RefCOCO+, and 5% in RefCOCOg,  
8 which undermines the authenticity of evaluations. We address this by excluding  
9 problematic instances and reevaluating several LMMs capable of handling the REC  
10 task, showing significant accuracy improvements, thus highlighting the impact  
11 of benchmark noise. In response, we introduce Ref-L4, a comprehensive REC  
12 benchmark, specifically designed to evaluate modern REC models. Ref-L4 is distinguished  
13 by four key features: 1) a substantial sample size with 45,341 annotations;  
14 2) a diverse range of object categories with 365 distinct types and varying instance  
15 scales from 30 to 3,767; 3) lengthy referring expressions averaging 24.2 words; and  
16 4) an extensive vocabulary comprising 22,813 unique words. We evaluate a total of  
17 24 large models on Ref-L4 and provide valuable insights. The cleaned versions of  
18 RefCOCO, RefCOCO+, and RefCOCOg, as well as our Ref-L4 benchmark and  
19 evaluation code, are available at <https://github.com/JierunChen/Ref-L4>.

## 20 1 Introduction

21 Referring expression comprehension (REC) [47, 38, 81, 21, 43, 75, 65] involves the task of localizing  
22 a specific target instance based on a given textual description. The advancement of REC has been  
23 significantly propelled by the superior language processing capabilities of large language models  
24 (LLMs) [55, 56, 37, 9, 15, 1, 28, 19, 26]. This progress is particularly evident in the exceptional  
25 performance of large multimodal models (LMMs) [62, 31, 13, 60, 2, 5, 66, 16, 57, 17, 80] on well-  
26 known benchmarks such as RefCOCO [71], RefCOCO+ [71], and RefCOCOg [36]. These models

---

\*Equal contribution.

†Corresponding author.

Table 1: Statistics of the labeling error rates for RefCOCO, RefCOCO+, and RefCOCOg, respectively. For each benchmark, the statistics are conducted on the combination of the validation and test sets.

Benchmark	Annotations	Errors	Labeling Error Rate
RefCOCO [71]	21,586	3,054	14%
RefCOCO+ [71]	21,373	5,201	24%
RefCOCOg [36]	14,498	675	5%

Table 2: The performance of four LMMs capable of handling the REC task on both the cleaned and original versions of the RefCOCO, RefCOCO+, and RefCOCOg benchmarks, using the conventional accuracy as the evaluation metric. The evaluation is performed on the combination of the validation and test sets for each benchmark. †: models fine-tuned on the specific dataset.

Benchmark	ONE-PEACE† [60]	OFA-L† [59]	OFA-L [59]	Qwen-VL [2]	CogVLM-Grounding [62]
RefCOCO [71]	92.15	89.85	85.13	88.51	92.44
RefCOCO (Cleaned)	94.11 (+1.96)	92.06 (+2.22)	87.95 (+2.81)	90.68 (+2.18)	94.58 (+2.13)
RefCOCO+ [71]	88.14	85.06	77.56	82.52	88.55
RefCOCO+ (Cleaned)	90.79 (+2.66)	87.38 (+2.32)	80.50 (+2.94)	85.60 (+3.08)	91.43 (+2.87)
RefCOCOg [36]	89.18	84.77	79.25	85.11	90.67
RefCOCOg (Cleaned)	90.75 (+1.57)	86.39 (+1.62)	80.89 (+1.64)	86.79 (+1.68)	92.36 (+1.68)

27 have demonstrated remarkable accuracy, with CogVLM [62], for instance, achieving an impressive  
28 accuracy rate of 92.44% on the RefCOCO benchmark.

29 This paper begins with a critical question: do existing REC benchmarks truly capture the comprehen-  
30 sive capabilities of LMMs? The foundational benchmarks, RefCOCO [71], RefCOCO+ [71], and  
31 RefCOCOg [36], were introduced sequentially in 2015, 2016, and 2016, respectively. In RefCOCO,  
32 the referring expressions are notably succinct, ranging from single words like “lady” and “yellow”  
33 to brief descriptions such as “far left person” and “white shirt”. RefCOCO+ intentionally excludes  
34 locational prepositions commonly found in RefCOCO, favoring short yet semantically rich expres-  
35 sions like “plastic cup with just ice” and “man on screen”. Conversely, RefCOCOg provides more  
36 elaborate annotations, including examples such as “a table of food, with plates, a pizza, pitchers, and  
37 glasses” and “a red and white checkered table with two wooden chairs”. These variations highlight the  
38 evolution and complexity of referring expressions across different benchmarks, raising the question  
39 of whether they can effectively assess the nuanced capabilities of modern LMMs in understanding  
40 diverse linguistic inputs and associating languages with visual elements.

41 **Labeling Error Rates of Existing Benchmarks.** To begin, we manually assess the labeling error  
42 rates of the validation and test sets in RefCOCO, RefCOCO+, and RefCOCOg, discovering a high  
43 error rate across these benchmarks. The labeling errors include, typos, misalignment between  
44 referring expressions and target instances, as well as inaccurate bounding box annotations, as depicted  
45 in Section A. As illustrated in Table 1, the labeling error rates for RefCOCO, RefCOCO+, and  
46 RefCOCOg are 14%, 24%, and 5%, respectively, indicating that evaluations performed on these  
47 benchmarks may lack authenticity.

48 **Reevaluation on RefCOCO, RefCOCO+ and RefCOCOg.** In response, we manually exclude  
49 the problematic instances from the validation and test sets of RefCOCO, RefCOCO+, and Ref-  
50 COCOg. Subsequently, we reevaluate four LMMs capable of handling the REC task—namely  
51 ONE-PEACE [60], OFA-L [59], Qwen-VL [2], and CogVLM-Grounding [62]—on both the cleaned  
52 and original versions of these datasets, as shown in Table 2. Across all models and cleaned bench-  
53 marks, we observe a significant accuracy improvement, ranging from 1.57 to 3.08, compared to their  
54 performance on the original versions. This demonstrates that noise in the benchmarks has impacted  
55 the models’ true capabilities. *To support further research in the REC field, we release the cleaned*  
56 *versions of RefCOCO, RefCOCO+, and RefCOCOg.*

57 **Ref-L4: A Comprehensive REC Benchmark for Modern LMM Evaluation.** We present Ref-L4,  
58 where L4 signifies four key aspects: a Large number of testing samples, Large diversity in object  
59 categories and instance scales, Long referring expressions, and a Large vocabulary. These features

Table 3: Comparison between our Ref-L4 benchmark and other REC benchmarks, including RefCOCO [71], RefCOCO+ [71], and RefCOCOg [36]. For the latter three benchmarks, we combine their validation and test sets for statistics. The instance size and image size are represented by their respective square roots. Avg. length: average length of annotations. Vocab.: vocabulary size.

Benchmark	Images	Instances	Annotations	Categories	Avg. Length	Instance Size	Image Size	Vocab.
RefCOCO [71]	3,000	7,596	21,586	71	3.6	105 - 607	230 - 640	3,525
RefCOCO+ [71]	3,000	7,578	21,373	71	3.6	105 - 607	230 - 640	4,387
RefCOCOg [36]	3,900	7,596	14,498	78	8.4	83 - 610	277 - 640	5,050
Ref-L4 (Ours)	9,735	18,653	45,341	365	24.2	30 - 3,767	230 - 6,606	22,813

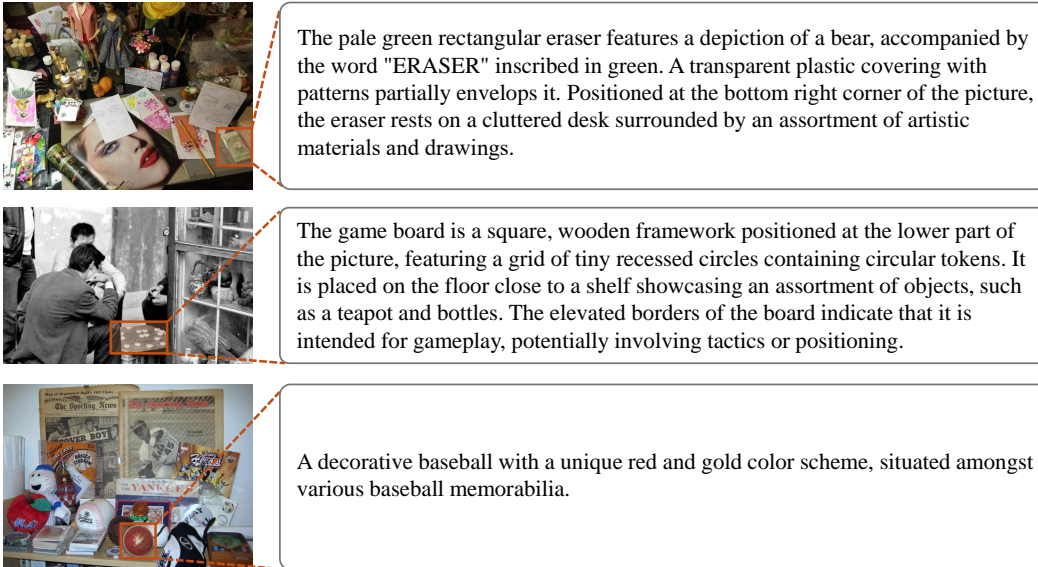


Figure 1: Examples from our Ref-L4 benchmark. We offer a detailed referring expression for each target instance represented by a bounding box. Zoom in for better visualization.

60 make Ref-L4 a comprehensive benchmark for assessing the REC capabilities of contemporary LMMs.  
 61 Table 3 provides a detailed comparison between Ref-L4 and other benchmarks including RefCOCO,  
 62 RefCOCO+, and RefCOCOg. Our Ref-L4 benchmark stands out due to the following characteristics:

- 63 • *Large-Scale.* Ref-L4 includes 9, 735 images, 18, 653 unique instances, and a total of 45, 341  
 64 annotations, significantly surpassing RefCOCO, RefCOCO+, and RefCOCOg. For instance,  
 65 RefCOCOg offers 3, 900 images, 7, 596 instances, and 14, 498 annotations.
- 66 • *High Diversity.* Ref-L4 features 365 unique categories. Since the RefCOCO series derive  
 67 from the COCO 2014 dataset, they encompass up to 78 categories. Additionally, our  
 68 benchmark covers a wider range of instance scales, from 30 to 3, 767, measured by the  
 69 square root of the instance area.
- 70 • *Lengthy Referring Expressions.* Each referring expression in Ref-L4 is a detailed description  
 71 of a specific instance, with lengths ranging from 3 to 117 words and an average of 24.2  
 72 words. In comparison, the average annotation lengths in RefCOCO, RefCOCO+, and  
 73 RefCOCOg are 3.6, 3.6, and 8.4 words, respectively. Examples can be found in Figure 1.
- 74 • *Extensive Vocabulary.* Due to the detailed nature of the referring expressions, Ref-L4  
 75 boasts a large vocabulary of 22, 813 words, which is four to six times larger than those of  
 76 RefCOCO, RefCOCO+, and RefCOCOg.

78 **Evaluation on Ref-L4.** We conduct an evaluation of 24 representative LMMs that can perform the  
 79 REC task. In addition to the standard accuracy metric, which considers predictions with an IoU  
 80 greater than 0.5 as accurate ( $\text{Acc}_{0.5}$ ), we also report accuracies at higher IoU thresholds:  $\text{Acc}_{0.75}$  and  
 81  $\text{Acc}_{0.9}$ . Furthermore, we introduce a mean accuracy (mAcc), calculated as the average accuracy from

82 Acc<sub>0.5</sub> to Acc<sub>0.9</sub> in increments of 0.05. To gain deeper insights into the models’ capabilities, we  
83 conduct a detailed analysis of REC performance across different instance scales and categories. *The*  
84 *Ref-L4 benchmark and the evaluation code are available at [https://github.com/JierunChen/](https://github.com/JierunChen/Ref-L4)*  
85 *Ref-L4*.

## 86 2 Related Work

87 **REC and Its Benchmarks.** Referring Expression Comprehension (REC) [47, 38, 81, 21, 43, 75] is a  
88 task that involves identifying a specific object within an image based on a given referring expression.  
89 Unlike object detection [30, 23, 52, 50, 4], which operates within fixed categories and a single  
90 visual modality, REC necessitates understanding free-form text to locate objects of any category.  
91 Phrase Grounding [44, 67, 14, 34, 27, 76, 61] is similar but typically involves shorter phrases and  
92 identifies multiple regions, whereas REC requires parsing longer expressions to pinpoint a single  
93 unique region. This complexity makes REC an ideal task for evaluating emerging large multimodal  
94 models. Current REC benchmarks such as RefCOCO [71], RefCOCO+[71], and RefCOCOg[36]  
95 include tens of thousands of annotations but are limited by their short expression lengths—averaging  
96 3.6, 3.6, and 8.4 words, respectively. Additionally, they encompass fewer than 80 categories, lacking  
97 real-world diversity. Other REC benchmarks [33, 8, 48, 7, 64, 24, 58, 10, 3, 12, 11, 18] are often  
98 designed for specific scenarios. For example, CLEVR-Ref+[33] focuses on simple objects like  
99 boxes, spheres, and cylinders. SK-VG[8] integrates prior scene knowledge as additional input, while  
100 RefCrowd [48] targets identifying a person within a crowd. By contrast, we introduce Ref-L4, a more  
101 general and comprehensive benchmark encompassing 365 categories and 45,341 annotations. Ref-L4  
102 features expressions averaging 24.2 words and a vocabulary of 22,813 words, facilitating the accurate  
103 evaluation of REC models on complex expressions and diverse objects.

104 **REC Models.** The evolution of REC models has transitioned from specialized models [20, 72,  
105 32, 54, 82, 68, 83] to generalist models or large multimodal models (LMMs)[62, 31, 13, 60, 2, 5,  
106 66, 78, 73, 74, 45, 77, 63, 53, 35, 46, 22]. Notable examples of these LMMs include CogVLM-  
107 Grounding[62], SPHINX [31, 13], ONE-PEACE [60], Qwen-VL-Chat [2], MiniGPTv2 [5], and  
108 Lenna [66]. These models, benefiting from larger model sizes and extensive training on diverse  
109 datasets, exhibit remarkable performance on conventional REC datasets. For example, CogVLM-  
110 Grounding achieves an accuracy of 94.58% on RefCOCO (cleaned). Additionally, the performance  
111 gap among models is shrinking, with many LMMs surpassing 90% accuracy. This performance  
112 saturation raises concerns about the adequacy of current REC benchmarks for making meaningful  
113 comparisons. In response, we propose Ref-L4, a more comprehensive and challenging benchmark.  
114 We have also conducted rigorous evaluations of 24 LMM models, offering holistic comparisons that  
115 highlight their weaknesses and suggest directions for improvement.

## 116 3 Ref-L4

### 117 3.1 Benchmark Creation

118 **Data Sources.** Our benchmark is derived from two sources: 1) our cleaned validation and test sets  
119 of the RefCOCO [71], RefCOCO+ [71], and RefCOCOg [36] datasets; and 2) the test set from the  
120 large-scale object detection dataset Objects365 [52]. The Objects365 dataset provides a broader  
121 range of categories, varying instance sizes, higher image resolutions, and more intricate scenes. In  
122 the RefCOCO series, each instance includes a bounding box, a category name, and an extremely brief  
123 expression like “right teddy bear”. In contrast, the Objects365 benchmark labels each instance with  
124 mainly a bounding box and the relevant category.

125 For the RefCOCO (cleaned) series, we begin by consolidating duplicate images and instances,  
126 resulting in a subset of 6,502 images containing 14,186 unique instances. For Objects365, we  
127 select samples from its testing set based on several criteria: 1) each image has both height  
128 and width greater than 800 pixels; 2) each image is sufficiently complex, containing more than  
129 10 categories and 20 instances; 3) each instance has a square normalized size  $\sqrt{(hw)/(HW)}$   
130 greater than 0.05, where  $(h, w)$  represents the instance size and  $(H, W)$  denotes the image size;  
131 4) we randomly sample  $N$  instances for each of the 365 classes defined in Objects365, with

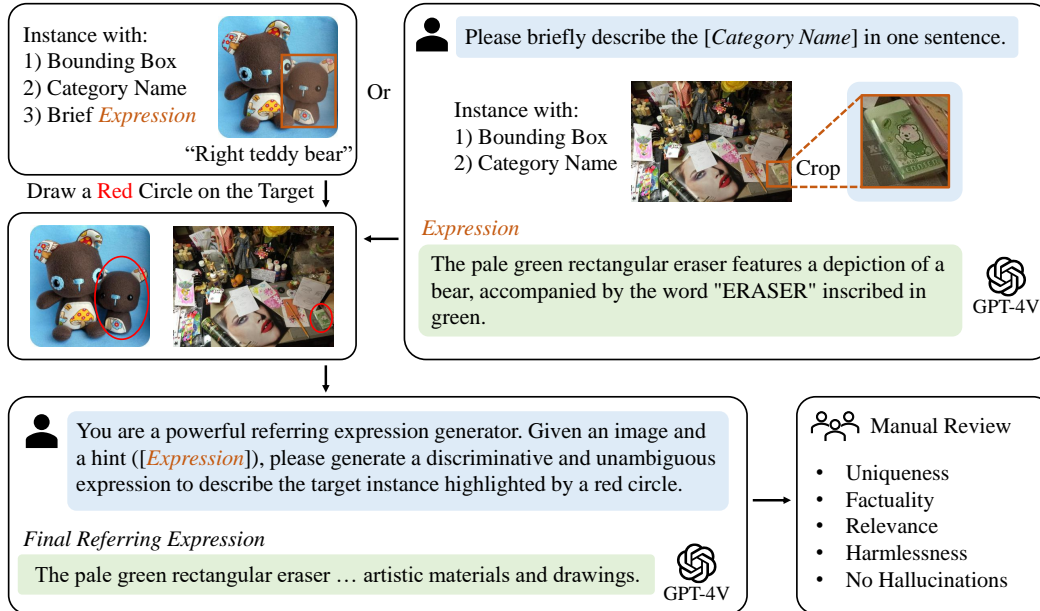


Figure 2: Pipeline of generating a referring expression for a target instance.

132  $N = \min(35, \text{the number of instances for the specific class}; 5)$  we review and exclude instances  
 133 with erroneous bounding box annotations or those difficult to describe uniquely. For a few rare  
 134 classes, we relax criterion-1 to 512 pixels and criterion-2 to 10 instances. Consequently, we collect  
 135 3, 233 images and 4, 467 instances from Objects365. Overall, our Ref-L4 benchmark comprises  
 136 9, 735 images and 18, 653 instances, sourced from the RefCOCO series and Objects365.

137 **Referring Expression Generation.** Given a target instance and its corresponding image, we leverage  
 138 GPT-4V with human reviewers in the loop to generate its precise and detailed referring expressions.  
 139 Figure 2 illustrates the three-step generation process:

140 *Step-1:* Each instance in the Objects365 dataset is linked to a bounding box and a *category name*. We  
 141 begin by cropping these instances from the original images. Next, we input each cropped area along  
 142 with the prompt detailed in Section B.1 into GPT-4V to produce a context-independent description.  
 143 For instances from the RefCOCO series, this step is omitted as each instance already has a brief  
 144 expression.

145 *Step-2:* Drawing inspiration from recent studies on GPT-4V [69], where GPT-4V is able to pay more  
 146 attention to instances highlighted by a red circle within an image, we similarly encircle the target  
 147 instance in red to facilitate GPT-4V in generating a context-aware referring expression. Following  
 148 this, as depicted in Figure 2, we process the image and use the prompt outlined in Section B.2 to  
 149 generate a context-aware referring expression for each instance. We instruct GPT-4V to describe  
 150 various features such as color, size, position, and context. Additionally, we provide a hint (the  
 151 context-independent description from Step-1) in the prompt to mitigate hallucination issues, resulting  
 152 in more accurate descriptions.

153 *Step-3:* We manually review all generated referring expressions to correct any hallucination issues.  
 154 We ensure that each expression uniquely describes the instance and is factual, accurate, and harmless.

155 **Annotation Expansion.** To date, we have compiled 18,653 unique referring expressions, each  
 156 describing a distinct instance. To assess the robustness of REC models to diverse language inputs, we  
 157 employ a two-stage rephrasing process to expand our benchmark: 1) utilizing GPT-4 with the prompt  
 158 detailed in Section B.3, to generate rephrased versions of each expression; 2) conducting a manual  
 159 review to ensure that the rephrased expressions are unique, factual, relevant, and harmless. Conse-  
 160 quently, our final Ref-L4 benchmark encompasses 9,735 images with 45,341 referring expressions,  
 161 each accurately describing one of the 18,653 unique instances.



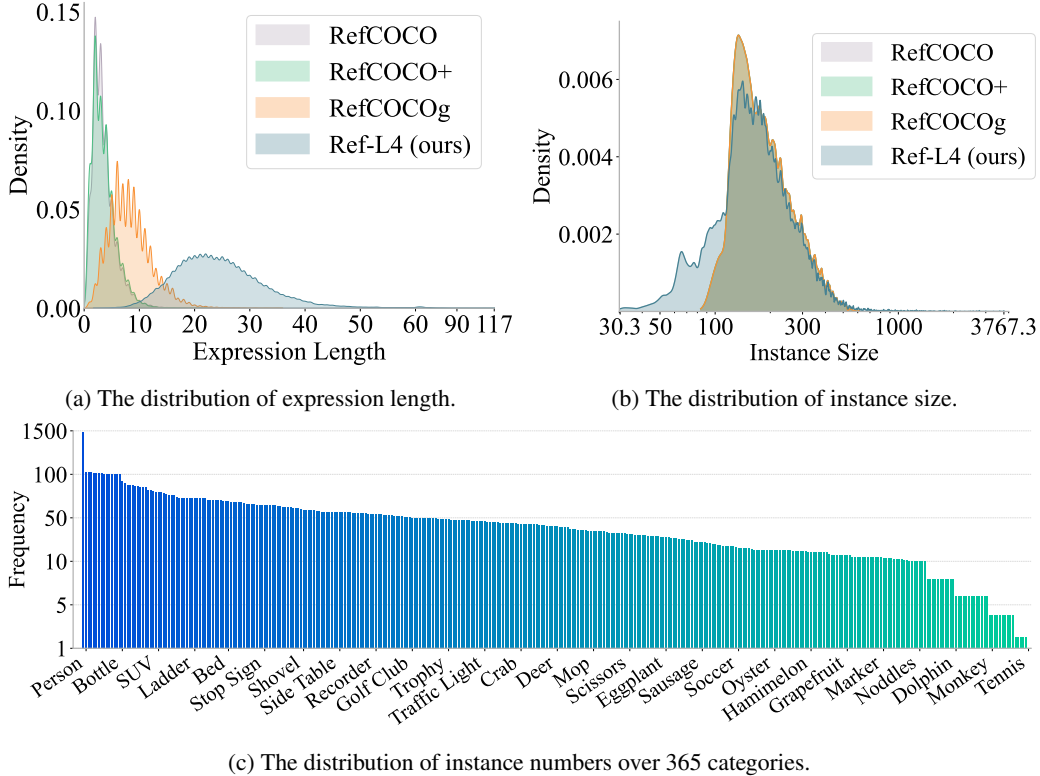


Figure 3: Analysis of referring expression length, instance size, and category distribution.

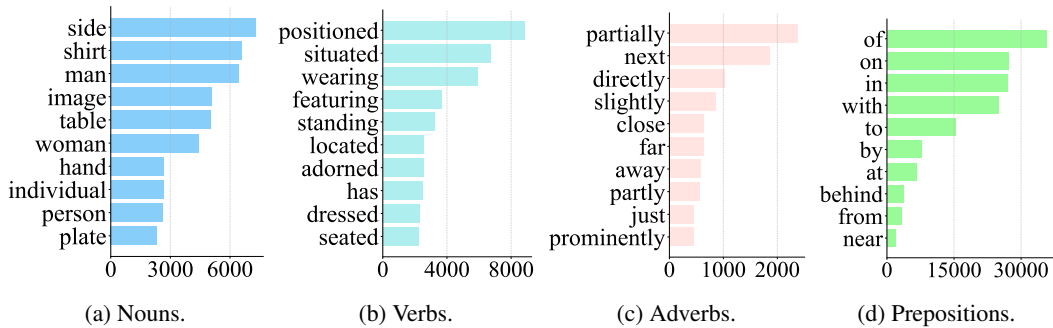


Figure 4: The frequency of the 10 most frequently used words in each part-of-speech category, as parsed using the SpaCy library.

### 162 3.2 Analysis

163 **Expression Length.** Figure 3a illustrates the distribution of expression lengths across four different  
 164 datasets: RefCOCO, RefCOCO+, RefCOCOg, and our Ref-L4. Due to the high overlap of data  
 165 samples, RefCOCO and RefCOCO+ exhibit similar distributions, with a high density of shorter  
 166 expressions peaking at around 3.6 words. RefCOCOg features slightly longer expressions on average,  
 167 peaking at approximately 8.4 words. In contrast, our Ref-L4 displays a significantly different  
 168 distribution, with expressions ranging much longer, peaking at around 24.2 words and having a long  
 169 tail extending up to 117 words. This suggests that our Ref-L4 benchmark is designed to push the  
 170 boundaries of current REC models, requiring them to process and comprehend more intricate and  
 171 detailed descriptions.

172 **Instance Size.** In Figure 3b, we present a density plot comparing the instance sizes across four  
 173 benchmarks. We define the instance size as the square root of the normalized size,  $\sqrt{(hw)/(HW)}$ ,  
 174 where  $(h, w)$  represents the dimensions of the instance and  $(H, W)$  represents the dimensions of the  
 175 image. All benchmarks exhibit a peak density around an instance size of 160. Our Ref-L4 benchmark

176 shows a wider distribution range compared to the other three, indicating that our Ref-L4 captures a  
177 broader spectrum of instance sizes.

178 **Categories.** Our Ref-4L benchmark comprises 18,653 instances spanning 365 distinct categories,  
179 providing more complex and diverse evaluation scenarios. In contrast, RefCOCO and RefCOCO+  
180 consists of 71 categories, while RefCOCOg covers 78 categories. Figure 3c presents the distribution  
181 of instances among these 365 categories. Notably, the ten categories with the highest number  
182 of instances are “Person”, “Chair”, “Hat”, “Desk”, “Lamp”, “Cabinet/shelf”, “Car”, “Sneakers”,  
183 “Handbag/Satchel”, and “Flag”.

184 **Vocabulary.** Our benchmark’s referring expressions comprise a vocabulary totaling 22,813 unique  
185 words. This is significantly larger than the vocabulary sizes of RefCOCO, RefCOCO+, and Ref-  
186 COCOg, which are 3,525, 4,387, and 5,050 words, respectively. Figure 4 illustrates the 10 most  
187 frequently used nouns, verbs, adverbs, and prepositions.

### 188 3.3 Evaluation

189 **Evaluation Metrics.** We propose three distinct evaluation protocols:

- 190 1. *Accuracy.* This is the conventional metric used in REC. For a given referring expression and  
191 corresponding image, the target instance is considered successfully localized if the IoU between  
192 the predicted bounding box and the ground truth exceeds 0.5. Accuracy is then calculated as the  
193 ratio of successfully localized samples to the total number of samples, referred to as  $\text{Acc}_{0.5}$  in  
194 this work. To better assess the localization capabilities of modern REC models, we also report  
195 accuracies at higher IoU thresholds:  $\text{Acc}_{0.75}$ ,  $\text{Acc}_{0.9}$ , and  $\text{mAcc}$ , which is the average accuracy  
196 from  $\text{Acc}_{0.5}$  to  $\text{Acc}_{0.9}$  in increments of 0.05.
- 197 2. *Scale-Aware Performance.* To gain deeper insights into model capabilities, we report performance  
198 based on instance sizes: small, medium, and large. The size of an instance is defined as the square  
199 root of its area,  $\sqrt{(hw)}$ , where  $(h, w)$  are the dimensions of the instance. Small instances are  
200 those with a size less than 128, medium instances are between 128 and 256, and large instances  
201 exceed 256. In total, there are 9345, 23280, and 12716 referring expressions describing 2, 954  
202 small, 10, 442 medium, and 5, 257 large instances, respectively.
- 203 3. *Per-Category Performance.* Our benchmark encompasses a wide range of categories, up to 365 in  
204 total. We provide an evaluation protocol to assess performance on a per-category basis.

205 **Benchmark Division.** Modern large multimodal models (LMMs) that are able to handle the REC  
206 task typically use unrestricted and extensive data for training. Our Ref-L4 benchmark is designed to  
207 assess the capabilities of these advanced models without imposing any limitations on the training data  
208 sources. The benchmark is divided into two subsets: a validation set, comprising 30% of the data  
209 with 7, 231 images, 10, 311 instances, and 13, 420 referring expressions; and a test set, comprising  
210 70% of the data with 9, 467 images, 17, 242 instances, and 31, 921 referring expressions. Given that  
211 our benchmark includes instances from 365 categories, we ensure that each category has at least one  
212 sample in both the validation and test sets. While we provide these two splits, we encourage the  
213 combined use of both sets for model evaluation, especially in the current LMM era, where the use of  
214 unrestricted training data is prevalent.

## 215 4 Experiments

216 **Main Result.** We evaluate a total of 24 LMMs that can perform the REC task, dividing them into  
217 two categories based on their output type: those that produce bounding boxes and those that produce  
218 segmentation masks. For models that output segmentation masks, we convert these masks into tight  
219 bounding boxes to enable evaluation on our Ref-L4 benchmark. Table 4 presents the performance  
220 of these models on the validation set, test set, and the combined set, using the metrics defined in  
221 Section 3.3. The evaluation prompt of GPT-4V is available in Section B.4. Among the models that  
222 output bounding boxes, CogVLM-Grounding [62] shows the best performance, while GlaMM [49]  
223 leads in performance among the models that output masks.

224 **Category-Wise Performance.** Each instance in our benchmark is assigned a category label from one  
225 of 365 classes. Figure 5 illustrates the performance of the top four models across these categories,

Table 4: Performance evaluation across 24 models on our Ref-L4 benchmark. NVIDIA A100 GPUs (80G) are utilized. The symbol † denotes models that outputs segmentation masks.

Model	Val+Test				Val	Test
	Acc <sub>0.5</sub>	Acc <sub>0.75</sub>	Acc <sub>0.9</sub>	mAcc	mAcc	mAcc
GPT-4V [39–41]	9.91	1.19	0.12	2.88	2.96	2.85
KOSMOS-2 [42]	48.53	38.34	17.54	34.72	34.89	34.64
OFA-Tiny [59]	55.21	43.22	27.70	41.44	41.53	41.40
OFA-Large [59]	72.53	62.31	45.02	59.17	59.42	59.07
Ferret-7b [70]	57.54	42.44	21.01	40.29	40.31	40.28
Ferret-13b [70]	64.44	49.04	27.46	46.88	47.31	46.71
GroundingGPT [29]	60.84	40.48	12.00	38.19	38.42	38.09
Shikra-7b [6]	65.06	39.62	10.45	38.60	38.91	38.47
Lenna [66]	65.90	58.55	45.58	55.69	55.88	55.60
MiniGPTv2 [5]	66.93	50.50	25.30	47.15	47.43	47.03
Qwen-VL-Chat [2]	73.80	58.05	37.16	55.94	56.18	55.83
ONE-PEACE [60]	70.82	60.09	36.12	55.07	55.49	54.89
SPHINX-MoE [13]	66.23	44.90	15.32	42.38	42.80	42.21
SPHINX-MoE-1k [13]	74.45	62.70	38.85	58.07	58.35	57.95
SPHINX [31]	74.78	53.65	21.15	50.09	50.33	49.99
SPHINX-1k [31]	78.52	62.17	32.95	57.57	57.91	57.42
SPHINX-v2-1k [31]	81.31	70.49	46.59	65.39	65.67	65.27
CogVLM-Grounding [62]	<b>81.70</b>	<b>70.77</b>	<b>48.35</b>	<b>66.09</b>	<b>66.25</b>	<b>66.02</b>
PixelLM-7B† [51]	41.83	27.57	13.32	27.10	27.09	27.11
PixelLM-13B† [51]	49.89	35.37	18.42	34.10	34.52	33.92
LISA-Explanatory† [25]	65.12	52.35	38.26	50.77	50.89	50.72
LISA† [25]	66.23	54.02	39.73	52.18	52.44	52.07
PSALM† [79]	67.26	58.22	44.11	55.46	55.68	55.37
GlaMM† [49]	<b>71.90</b>	<b>60.27</b>	<b>45.15</b>	<b>57.89</b>	<b>58.16</b>	<b>57.78</b>

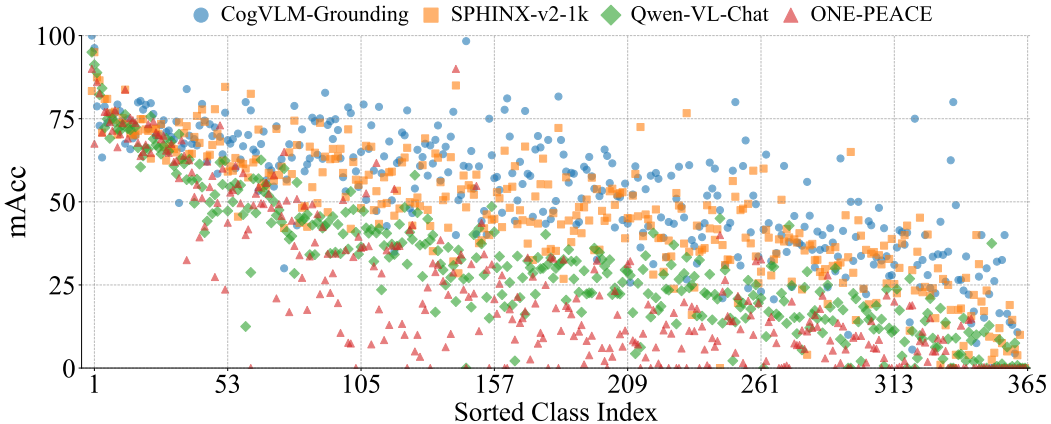


Figure 5: Category-wise performance of the four top-performing models on the val+test set, sorted in descending order based on their average per-category performance. The performance of all models can be found in Section C.1.

226 sorted in descending order based on their average per-category performance. The results indicate a  
 227 training bias issue, as all four models exhibit poor performance on some common categories.

228 **Scale-Aware Evaluation.** In Section 3.3, we present a scale-aware evaluation to assess the model’s  
 229 ability to handle different instance scales. Specifically, we categorize all samples in our benchmark  
 230 into three sets based on instance size: small, medium, and large. The performance of 24 models is  
 231 detailed in Table 5. Among the bounding-box-output models, CogVLM-Grounding [62] excels with  
 232 small and medium instances, while SPHINX-v2-1k [31] achieves the best performance with large  
 233 instances. For mask-output models, GlaMM [49] outperforms all other models across all three sets.

234 **Evaluation on Diverse Data Sources.** Our benchmark is derived from COCO and Objects365  
 235 datasets. We assess the performance of the top four models with bounding box outputs and the top



Table 5: Scale-aware evaluation across 24 models on our Ref-L4 benchmark.

Model	Small Size		Medium Size		Large Size	
	Acc <sub>0.5</sub>	mAcc	Acc <sub>0.5</sub>	mAcc	Acc <sub>0.5</sub>	mAcc
GPT-4V [39–41]	2.13	0.49	10.29	2.78	14.93	4.83
KOSMOS-2 [42]	24.19	11.63	46.95	32.91	69.32	54.98
OFA-Tiny [59]	17.91	11.49	65.13	49.00	64.46	49.61
OFA-Large [59]	40.13	27.07	81.03	66.49	80.78	69.36
Ferret-7b [70]	30.93	14.57	62.40	43.72	68.18	52.92
Ferret-13b [70]	36.46	17.88	70.50	51.86	73.92	59.09
GroundingGPT [29]	24.43	10.28	67.67	41.04	75.09	53.47
Shikra-7b [6]	43.91	18.50	75.98	46.27	60.60	39.34
Lenna [66]	31.02	23.48	72.90	61.53	78.72	68.66
MiniGPTv2 [5]	32.99	14.85	73.67	51.16	79.52	63.53
Qwen-VL-Chat [2]	47.66	26.26	79.80	61.06	82.01	68.37
ONE-PEACE [60]	22.18	13.98	83.26	63.39	83.81	70.04
SPHINX-MoE [13]	39.48	16.39	72.97	46.38	73.55	54.17
SPHINX-MoE-1k [13]	58.96	37.61	77.80	61.53	79.70	66.77
SPHINX [31]	48.82	22.08	80.56	54.10	83.27	63.34
SPHINX-1k [31]	59.48	33.21	82.95	61.82	84.40	67.68
SPHINX-v2-1k [31]	65.23	43.43	84.00	68.45	<b>88.21</b>	<b>75.91</b>
CogVLM-Grounding [62]	<b>75.06</b>	<b>52.85</b>	<b>86.43</b>	<b>71.31</b>	77.91	66.25
PixelLM-7B <sup>†</sup> [51]	8.25	4.05	43.90	27.33	62.72	43.64
PixelLM-13B <sup>†</sup> [51]	17.05	8.54	53.40	35.48	67.59	50.34
LISA-Explanatory <sup>†</sup> [25]	39.11	27.16	70.03	54.61	75.25	61.09
LISA <sup>†</sup> [25]	39.24	27.49	71.17	56.05	77.01	63.22
PSALM <sup>†</sup> [79]	37.35	28.43	75.06	61.79	74.97	63.74
GlaMM <sup>†</sup> [49]	<b>47.07</b>	<b>34.36</b>	<b>77.17</b>	<b>62.28</b>	<b>80.50</b>	<b>67.14</b>

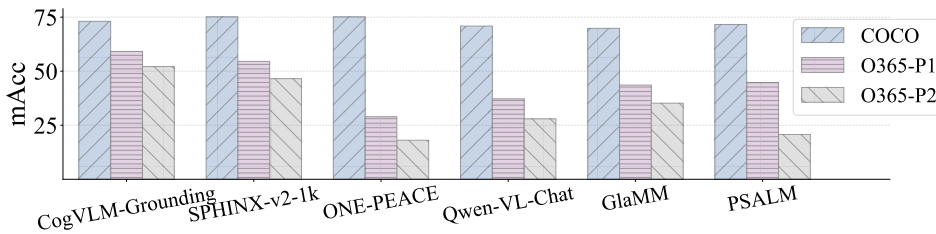


Figure 6: Evaluation of six models on various data sources, with mAcc acting as the metric. The results of all models can be found in Section C.2.

236 two models with mask outputs across various subsets originating from either COCO or Objects365.  
 237 These subsets are: 1) the COCO-derived set (referred to as “COCO”); 2) a subset from Objects365,  
 238 where the instances have categories that also exist in COCO (referred to as “O365-P1”); 3) another  
 239 subset from Objects365, where the instances have categories not found in COCO (referred to as  
 240 “O365-P2”). Figure 6 presents the performance of these models across the three subsets. The “COCO”  
 241 set shows higher accuracy compared to the other two sets, partially because most models are trained  
 242 on the RefCOCO series and have limited exposure to Objects365 images. “O365-P1” exhibits higher  
 243 accuracy than “O365-P2”, as the latter includes more rare categories.

## 244 5 Conclusion

245 In this work, we first point out several limitations of the current REC benchmarks, such as substantial  
 246 labeling inaccuracies and very brief referring expressions. To better assess the capabilities of models,  
 247 particularly those LMMs that can perform the REC task, we present Ref-L4, which features four key  
 248 characteristics: 1) a large-scale dataset with 45,341 annotations; 2) a wide range of object categories  
 249 and varying instance scales; 3) detailed referring expressions; and 4) an extensive vocabulary  
 250 comprising 22,813 unique words. We evaluate a total of 24 models using various evaluation protocols.  
 251 We wish that Ref-L4 could serve as a valuable resource for researchers and developers, fostering the  
 252 development of more robust and versatile REC models in the LMM era.

## 253 References

- 254 [1] J. Bai, S. Bai, Y. Chu, Z. Cui, K. Dang, X. Deng, Y. Fan, W. Ge, Y. Han, F. Huang, et al. Qwen technical  
255 report. *arXiv preprint arXiv:2309.16609*, 2023.
- 256 [2] J. Bai, S. Bai, S. Yang, S. Wang, S. Tan, P. Wang, J. Lin, C. Zhou, and J. Zhou. Qwen-vl: A frontier large  
257 vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023.
- 258 [3] Y. Bu, L. Li, J. Xie, Q. Liu, Y. Cai, Q. Huang, and Q. Li. Scene-text oriented referring expression  
259 comprehension. *IEEE Transactions on Multimedia*, 2022.
- 260 [4] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko. End-to-end object detection  
261 with transformers. In *European conference on computer vision*, pages 213–229. Springer, 2020.
- 262 [5] J. Chen, D. Zhu, X. Shen, X. Li, Z. Liu, P. Zhang, R. Krishnamoorthi, V. Chandra, Y. Xiong, and  
263 M. Elhoseiny. Minigt-v2: large language model as a unified interface for vision-language multi-task  
264 learning. *arXiv preprint arXiv:2310.09478*, 2023.
- 265 [6] K. Chen, Z. Zhang, W. Zeng, R. Zhang, F. Zhu, and R. Zhao. Shikra: Unleashing multimodal llm’s  
266 referential dialogue magic. *arXiv preprint arXiv:2306.15195*, 2023.
- 267 [7] Z. Chen, P. Wang, L. Ma, K.-Y. K. Wong, and Q. Wu. Cops-ref: A new dataset and task on compositional  
268 referring expression comprehension. In *Proceedings of the IEEE/CVF Conference on Computer Vision  
269 and Pattern Recognition*, pages 10086–10095, 2020.
- 270 [8] Z. Chen, R. Zhang, Y. Song, X. Wan, and G. Li. Advancing visual grounding with scene knowledge:  
271 Benchmark and method. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
272 Recognition*, pages 15039–15049, 2023.
- 273 [9] W.-L. Chiang, Z. Li, Z. Lin, Y. Sheng, Z. Wu, H. Zhang, L. Zheng, S. Zhuang, Y. Zhuang, J. E. Gonzalez,  
274 et al. Vicuna: An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality. *See [https://vicuna.  
275 lmsys.org](https://vicuna.lmsys.org) (accessed 14 April 2023)*, 2(3):6, 2023.
- 276 [10] V. Cirik, T. Berg-Kirkpatrick, and L.-P. Morency. Refer360 degree: A referring expression recogni-  
277 tion dataset in 360 degree images. In *Proceedings of the 58th Annual Meeting of the Association for  
278 Computational Linguistics*, pages 7189–7202, 2020.
- 279 [11] H. De Vries, F. Strub, S. Chandar, O. Pietquin, H. Larochelle, and A. Courville. Guesswhat?! visual object  
280 discovery through multi-modal dialogue. In *Proceedings of the IEEE Conference on Computer Vision and  
281 Pattern Recognition*, pages 5503–5512, 2017.
- 282 [12] C. Gao, B. Yang, H. Wang, M. Yang, W. Yu, Y. Liu, and X. Bai. Textrec: A dataset for referring expression  
283 comprehension with reading comprehension. In *International Conference on Document Analysis and  
284 Recognition*, pages 402–420. Springer, 2023.
- 285 [13] P. Gao, R. Zhang, C. Liu, L. Qiu, S. Huang, W. Lin, S. Zhao, S. Geng, Z. Lin, P. Jin, et al. Sphinx-x: Scaling  
286 data and parameters for a family of multi-modal large language models. *arXiv preprint arXiv:2402.05935*,  
287 2024.
- 288 [14] A. Gupta, P. Dollar, and R. Girshick. Lvis: A dataset for large vocabulary instance segmentation. In  
289 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5356–5364,  
290 2019.
- 291 [15] Y. Hao, H. Song, L. Dong, S. Huang, Z. Chi, W. Wang, S. Ma, and F. Wei. Language models are  
292 general-purpose interfaces. *arXiv preprint arXiv:2206.06336*, 2022.
- 293 [16] J. He, Y. Wang, L. Wang, H. Lu, J.-Y. He, J.-P. Lan, B. Luo, and X. Xie. Multi-modal instruction tuned  
294 llms with fine-grained visual perception. *arXiv preprint arXiv:2403.02969*, 2024.
- 295 [17] Z. Huang, Z. Zhang, Z.-J. Zha, Y. Lu, and B. Guo. Relationvlm: Making large vision-language models  
296 understand visual relations. *arXiv preprint arXiv:2403.12801*, 2024.
- 297 [18] B. Jia, Y. Chen, H. Yu, Y. Wang, X. Niu, T. Liu, Q. Li, and S. Huang. Sceneverse: Scaling 3d vision-  
298 language learning for grounded scene understanding. *arXiv preprint arXiv:2401.09340*, 2024.
- 299 [19] A. Q. Jiang, A. Sablayrolles, A. Roux, A. Mensch, B. Savary, C. Bamford, D. S. Chaplot, D. d. l. Casas,  
300 E. B. Hanna, F. Bressand, et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.

- 301 [20] A. Kamath, M. Singh, Y. LeCun, G. Synnaeve, I. Misra, and N. Carion. Mdetr-modulated detection for  
302 end-to-end multi-modal understanding. In *Proceedings of the IEEE/CVF International Conference on*  
303 *Computer Vision*, pages 1780–1790, 2021.
- 304 [21] S. Kazemzadeh, V. Ordonez, M. Matten, and T. Berg. Referitgame: Referring to objects in photographs of  
305 natural scenes. In *Proceedings of the 2014 conference on empirical methods in natural language processing*  
306 *(EMNLP)*, pages 787–798, 2014.
- 307 [22] M. KOSAREVA. Pushing the limits of visual grounding: Pre-training on large synthetic datasets.
- 308 [23] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma,  
309 et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations.  
310 *International journal of computer vision*, 123:32–73, 2017.
- 311 [24] S. Kurita, N. Katsura, and E. Onami. Refego: Referring expression comprehension dataset from first-person  
312 perception of ego4d. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,  
313 pages 15214–15224, 2023.
- 314 [25] X. Lai, Z. Tian, Y. Chen, Y. Li, Y. Yuan, S. Liu, and J. Jia. Lisa: Reasoning segmentation via large language  
315 model. *arXiv preprint arXiv:2308.00692*, 2023.
- 316 [26] M. Lewis, Y. Liu, N. Goyal, M. Ghazvininejad, A. Mohamed, O. Levy, V. Stoyanov, and L. Zettlemoyer.  
317 Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and  
318 comprehension. *arXiv preprint arXiv:1910.13461*, 2019.
- 319 [27] L. H. Li, P. Zhang, H. Zhang, J. Yang, C. Li, Y. Zhong, L. Wang, L. Yuan, L. Zhang, J.-N. Hwang, et al.  
320 Grounded language-image pre-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision*  
321 *and Pattern Recognition*, pages 10965–10975, 2022.
- 322 [28] Y. Li, S. Bubeck, R. Eldan, A. Del Giorno, S. Gunasekar, and Y. T. Lee. Textbooks are all you need ii:  
323 phi-1.5 technical report. *arXiv preprint arXiv:2309.05463*, 2023.
- 324 [29] Z. Li, Q. Xu, D. Zhang, H. Song, Y. Cai, Q. Qi, R. Zhou, J. Pan, Z. Li, V. T. Vu, et al. Lego: Language  
325 enhanced multi-modal grounding model. *arXiv preprint arXiv:2401.06071*, 2024.
- 326 [30] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft  
327 coco: Common objects in context. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich,*  
328 *Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014.
- 329 [31] Z. Lin, C. Liu, R. Zhang, P. Gao, L. Qiu, H. Xiao, H. Qiu, C. Lin, W. Shao, K. Chen, et al. Sphinx: The  
330 joint mixing of weights, tasks, and visual embeddings for multi-modal large language models. *arXiv*  
331 *preprint arXiv:2311.07575*, 2023.
- 332 [32] J. Liu, L. Wang, and M.-H. Yang. Referring expression generation and comprehension via attributes. In  
333 *Proceedings of the IEEE International Conference on Computer Vision*, pages 4856–4864, 2017.
- 334 [33] R. Liu, C. Liu, Y. Bai, and A. L. Yuille. Clevr-ref+: Diagnosing visual reasoning with referring expressions.  
335 In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4185–4194,  
336 2019.
- 337 [34] S. Liu, Z. Zeng, T. Ren, F. Li, H. Zhang, J. Yang, C. Li, J. Yang, H. Su, J. Zhu, et al. Grounding dino:  
338 Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*,  
339 2023.
- 340 [35] C. Ma, Y. Jiang, J. Wu, Z. Yuan, and X. Qi. Groma: Localized visual tokenization for grounding multimodal  
341 large language models. *arXiv preprint arXiv:2404.13013*, 2024.
- 342 [36] J. Mao, J. Huang, A. Toshev, O. Camburu, A. L. Yuille, and K. Murphy. Generation and comprehension of  
343 unambiguous object descriptions. In *Proceedings of the IEEE conference on computer vision and pattern*  
344 *recognition*, pages 11–20, 2016.
- 345 [37] A. Meta. Introducing meta llama 3: The most capable openly available llm to date. *Meta AI*, 2024.
- 346 [38] V. K. Nagaraja, V. I. Morariu, and L. S. Davis. Modeling context between objects for referring expression  
347 understanding. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands,*  
348 *October 11–14, 2016, Proceedings, Part IV 14*, pages 792–807. Springer, 2016.
- 349 [39] OpenAI. Gpt-4 technical report, 2023.

- 350 [40] OpenAI. Gpt-4v(ision) system card. 2023. URL [https://cdn.openai.com/papers/GPTV\\_System\\_Card.pdf](https://cdn.openai.com/papers/GPTV_System_Card.pdf).  
351
- 352 [41] OpenAI. Gpt-4v(ision) technical work and authors. 2023. URL <https://cdn.openai.com/contributions/gpt-4v.pdf>.  
353
- 354 [42] Z. Peng, W. Wang, L. Dong, Y. Hao, S. Huang, S. Ma, and F. Wei. Kosmos-2: Grounding multimodal large language models to the world. *arXiv preprint arXiv:2306.14824*, 2023.  
355
- 356 [43] R. Pi, L. Yao, J. Gao, J. Zhang, and T. Zhang. Perceptiongpt: Effectively fusing visual perception into llm. *arXiv preprint arXiv:2311.06612*, 2023.  
357
- 358 [44] B. A. Plummer, L. Wang, C. M. Cervantes, J. C. Caicedo, J. Hockenmaier, and S. Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In *Proceedings of the IEEE international conference on computer vision*, pages 2641–2649, 2015.  
359  
360
- 361 [45] S. Pramanick, G. Han, R. Hou, S. Nag, S.-N. Lim, N. Ballas, Q. Wang, R. Chellappa, and A. Almahairi. Jack of all tasks, master of many: Designing general-purpose coarse-to-fine vision-language model. *arXiv preprint arXiv:2312.12423*, 2023.  
362  
363
- 364 [46] L. Qi, Y.-W. Chen, L. Yang, T. Shen, X. Li, W. Guo, Y. Xu, and M.-H. Yang. Generalizable entity grounding via assistance of large language model. *arXiv preprint arXiv:2402.02555*, 2024.  
365
- 366 [47] Y. Qiao, C. Deng, and Q. Wu. Referring expression comprehension: A survey of methods and datasets. *IEEE Transactions on Multimedia*, 23:4426–4440, 2020.  
367
- 368 [48] H. Qiu, H. Li, T. Zhao, L. Wang, Q. Wu, and F. Meng. Refcrowd: Grounding the target in crowd with referring expressions. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 4435–4444, 2022.  
369  
370
- 371 [49] H. Rasheed, M. Maaz, S. Shaji, A. Shaker, S. Khan, H. Cholakkal, R. M. Anwer, E. Xing, M.-H. Yang, and F. S. Khan. Glamm: Pixel grounding large multimodal model. *arXiv preprint arXiv:2311.03356*, 2023.  
372
- 373 [50] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016.  
374
- 375 [51] Z. Ren, Z. Huang, Y. Wei, Y. Zhao, D. Fu, J. Feng, and X. Jin. Pixellm: Pixel reasoning with large multimodal model. *arXiv preprint arXiv:2312.02228*, 2023.  
376
- 377 [52] S. Shao, Z. Li, T. Zhang, C. Peng, G. Yu, X. Zhang, J. Li, and J. 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.  
378  
379
- 380 [53] H. Shen, T. Zhao, M. Zhu, and J. Yin. Groundvlp: Harnessing zero-shot visual grounding from vision-language pre-training and open-vocabulary object detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 4766–4775, 2024.  
381  
382
- 383 [54] W. Su, X. Zhu, Y. Cao, B. Li, L. Lu, F. Wei, and J. Dai. Vi-bert: Pre-training of generic visual-linguistic representations. *arXiv preprint arXiv:1908.08530*, 2019.  
384
- 385 [55] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.  
386  
387
- 388 [56] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.  
389  
390
- 391 [57] H. Wang, H. Tang, L. Jiang, S. Shi, M. F. Naeem, H. Li, B. Schiele, and L. Wang. Git: Towards generalist vision transformer through universal language interface. *arXiv preprint arXiv:2403.09394*, 2024.  
392
- 393 [58] P. Wang, D. Liu, H. Li, and Q. Wu. Give me something to eat: referring expression comprehension with commonsense knowledge. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 28–36, 2020.  
394  
395
- 396 [59] P. Wang, A. Yang, R. Men, J. Lin, S. Bai, Z. Li, J. Ma, C. Zhou, J. Zhou, and H. Yang. Ofa: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. In *International Conference on Machine Learning*, pages 23318–23340. PMLR, 2022.  
397  
398

- 399 [60] P. Wang, S. Wang, J. Lin, S. Bai, X. Zhou, J. Zhou, X. Wang, and C. Zhou. One-peace: Exploring one  
400 general representation model toward unlimited modalities. *arXiv preprint arXiv:2305.11172*, 2023.
- 401 [61] Q. Wang, H. Tan, S. Shen, M. W. Mahoney, and Z. Yao. Maf: Multimodal alignment framework for  
402 weakly-supervised phrase grounding. *arXiv preprint arXiv:2010.05379*, 2020.
- 403 [62] W. Wang, Q. Lv, W. Yu, W. Hong, J. Qi, Y. Wang, J. Ji, Z. Yang, L. Zhao, X. Song, et al. Cogvlm: Visual  
404 expert for pretrained language models. *arXiv preprint arXiv:2311.03079*, 2023.
- 405 [63] W. Wang, Z. Chen, X. Chen, J. Wu, X. Zhu, G. Zeng, P. Luo, T. Lu, J. Zhou, Y. Qiao, et al. Visionllm: Large  
406 language model is also an open-ended decoder for vision-centric tasks. *Advances in Neural Information  
407 Processing Systems*, 36, 2024.
- 408 [64] W. Wang, Y. Zhang, X. He, Y. Yan, Z. Zhao, X. Wang, and J. Liu. Beyond literal descriptions: Understand-  
409 ing and locating open-world objects aligned with human intentions. *arXiv preprint arXiv:2402.11265*,  
410 2024.
- 411 [65] Y. Wang, Z. Ji, D. Wang, Y. Pang, and X. Li. Towards unsupervised referring expression comprehension  
412 with visual semantic parsing. *Knowledge-Based Systems*, 285:111318, 2024.
- 413 [66] F. Wei, X. Zhang, A. Zhang, B. Zhang, and X. Chu. Lenna: Language enhanced reasoning detection  
414 assistant. *arXiv preprint arXiv:2312.02433*, 2023.
- 415 [67] C. Wu, Z. Lin, S. Cohen, T. Bui, and S. Maji. Phrasecut: Language-based image segmentation in the  
416 wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages  
417 10216–10225, 2020.
- 418 [68] B. Yan, Y. Jiang, J. Wu, D. Wang, P. Luo, Z. Yuan, and H. Lu. Universal instance perception as object  
419 discovery and retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
420 Recognition*, pages 15325–15336, 2023.
- 421 [69] Z. Yang, L. Li, K. Lin, J. Wang, C.-C. Lin, Z. Liu, and L. Wang. The dawn of lmms: Preliminary  
422 explorations with gpt-4v (ision). *arXiv preprint arXiv:2309.17421*, 9(1):1, 2023.
- 423 [70] H. You, H. Zhang, Z. Gan, X. Du, B. Zhang, Z. Wang, L. Cao, S.-F. Chang, and Y. Yang. Ferret: Refer and  
424 ground anything anywhere at any granularity. *arXiv preprint arXiv:2310.07704*, 2023.
- 425 [71] L. Yu, P. Poirson, S. Yang, A. C. Berg, and T. L. Berg. Modeling context in referring expressions. In  
426 *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14,  
427 2016, Proceedings, Part II 14*, pages 69–85. Springer, 2016.
- 428 [72] L. Yu, Z. Lin, X. Shen, J. Yang, X. Lu, M. Bansal, and T. L. Berg. Mattnet: Modular attention network  
429 for referring expression comprehension. In *Proceedings of the IEEE conference on computer vision and  
430 pattern recognition*, pages 1307–1315, 2018.
- 431 [73] Y. Zhan, Y. Zhu, Z. Chen, F. Yang, M. Tang, and J. Wang. Griffon: Spelling out all object locations at any  
432 granularity with large language models. *arXiv preprint arXiv:2311.14552*, 2023.
- 433 [74] Y. Zhan, Y. Zhu, H. Zhao, F. Yang, M. Tang, and J. Wang. Griffon v2: Advancing multimodal perception  
434 with high-resolution scaling and visual-language co-referring. *arXiv preprint arXiv:2403.09333*, 2024.
- 435 [75] C. Zhang, W. Li, W. Ouyang, Q. Wang, W.-S. Kim, and S. Hong. Referring expression comprehension with  
436 semantic visual relationship and word mapping. In *Proceedings of the 27th ACM International Conference  
437 on Multimedia*, pages 1258–1266, 2019.
- 438 [76] H. Zhang, P. Zhang, X. Hu, Y.-C. Chen, L. Li, X. Dai, L. Wang, L. Yuan, J.-N. Hwang, and J. Gao. Glipv2:  
439 Unifying localization and vision-language understanding. *Advances in Neural Information Processing  
440 Systems*, 35:36067–36080, 2022.
- 441 [77] H. Zhang, H. Li, F. Li, T. Ren, X. Zou, S. Liu, S. Huang, J. Gao, L. Zhang, C. Li, et al. Llava-grounding:  
442 Grounded visual chat with large multimodal models. *arXiv preprint arXiv:2312.02949*, 2023.
- 443 [78] H. Zhang, H. You, P. Dufter, B. Zhang, C. Chen, H.-Y. Chen, T.-J. Fu, W. Y. Wang, S.-F. Chang, Z. Gan,  
444 et al. Ferret-v2: An improved baseline for referring and grounding with large language models. *arXiv  
445 preprint arXiv:2404.07973*, 2024.
- 446 [79] Z. Zhang, Y. Ma, E. Zhang, and X. Bai. Psalm: Pixelwise segmentation with large multi-modal model.  
447 *arXiv preprint arXiv:2403.14598*, 2024.



- 448 [80] H. Zhao, W. Ge, and Y.-c. Chen. Llm-optic: Unveiling the capabilities of large language models for  
449 universal visual grounding. *arXiv preprint arXiv:2405.17104*, 2024.
- 450 [81] D. Zheng, T. Kong, Y. Jing, J. Wang, and X. Wang. Towards unifying reference expression generation and  
451 comprehension. *arXiv preprint arXiv:2210.13076*, 2022.
- 452 [82] Z. Zheng, W. Wang, S. Qi, and S.-C. Zhu. Reasoning visual dialogs with structural and partial observations.  
453 In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6669–6678,  
454 2019.
- 455 [83] X. Zou, J. Yang, H. Zhang, F. Li, L. Li, J. Wang, L. Wang, J. Gao, and Y. J. Lee. Segment everything  
456 everywhere all at once. *Advances in Neural Information Processing Systems*, 36, 2024.

## 457 A Labeling Errors in Existing Benchmarks

458 In the REC task, a referring expression should uniquely describe an instance, which is represented by  
459 an accurate bounding box. We have identified and visualized three common types of labeling errors  
460 in the RefCOCO, RefCOCO+, and RefCOCOg benchmarks: 1) non-unique referring expressions  
461 (Figure 7), which refer to multiple instances within the same image; 2) inaccurate bounding boxes  
462 (Figure 8); and 3) misalignment between target instances and their referring expressions (Figure 9),  
463 where the referring expressions are either ambiguous or do not refer to any instance in the image.

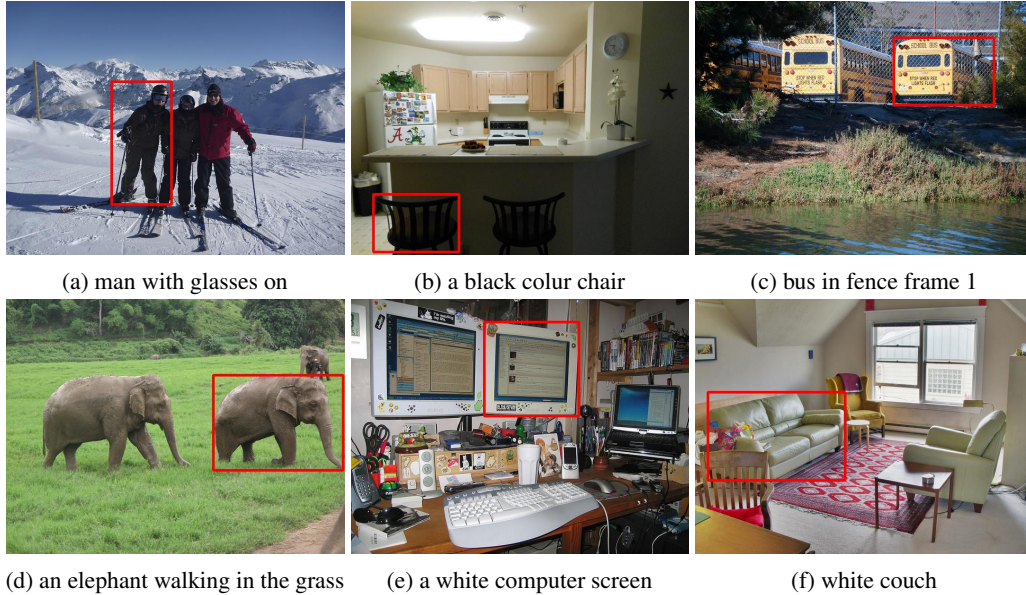


Figure 7: Visualization of labeling errors, where a referring expression refers to multiple instances within the same image. For each sub-figure, we display the original bounding box annotation with a red rectangle and include the corresponding referring expression in the caption.

## 464 B Prompts

### 465 B.1 Prompt for Context-Independent Description Generation

466 Briefly describe the [*Category Name*] in one sentence. Begin your description with the object name,  
467 including adjectives if appropriate to describe its color or shape. Focus only on visible features and  
468 avoid mentioning blurriness.

469 Input image: [*Cropped Image*].

### 470 B.2 Prompt for Context-Aware Description Generation

471 You are a sophisticated referring expression generator. Your task is to generate a clear and specific  
472 description for the target instance highlighted by a red circle in the provided image, based on a given  
473 hint and the following criteria:

474 *Criteria 1:* The description should enable individuals to understand and accurately identify the  
475 specified region within the image.

476 *Criteria 2:* The description may should various attributes such as category, shape, size, color,  
477 visibility, exposure, texture, orientation, absolute position, relative position, facial features, clothing,  
478 accessories, gestures, context, semantic attributes, emotions, age, gender, posture, action, and  
479 especially interactions with other instances. The selection of features should be relevant to the  
480 particular region and the image context.

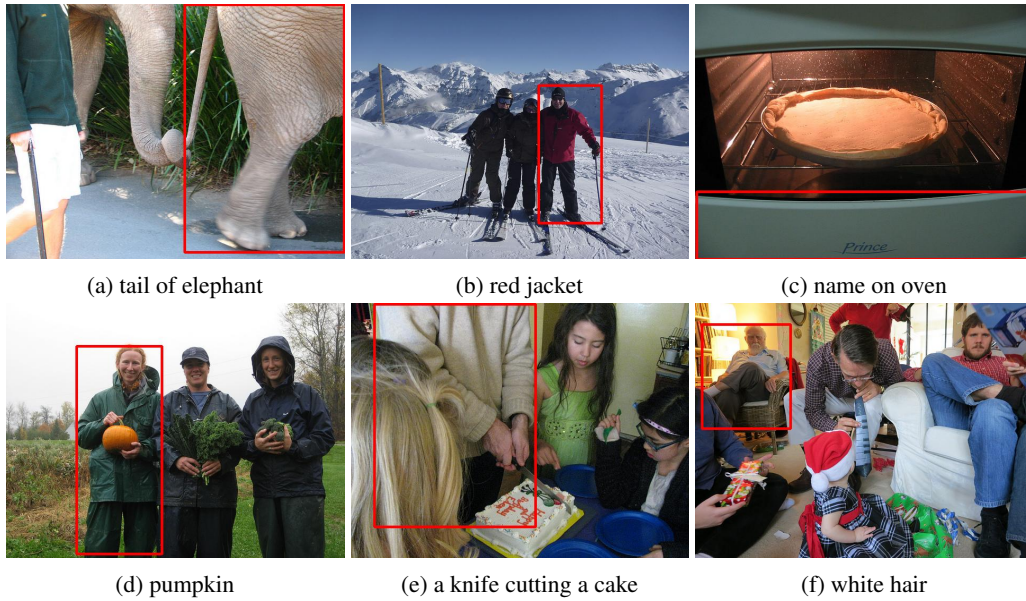


Figure 8: Visualization of labeling errors, where the bounding box annotations are inaccurate. For each sub-figure, we display the original bounding box annotation with a red rectangle and include the corresponding referring expression in the caption.



Figure 9: Visualization of labeling errors, where the referring expressions are either ambiguous or do not refer to any instance in the image. For each sub-figure, we display the original bounding box annotation with a red rectangle and include the corresponding referring expression in the caption.

481 *Criteria 3:* The red circle is solely for highlighting the region of interest. Do not refer to it in your  
 482 descriptions.

483 *Criteria 4:* Avoid using unnecessary words like “look for”, “spot”, “observe”, “find”, “notice”,  
 484 “identify”, “outline”, “target” and “question”.

485 *Criteria 5:* Ensure that the subject of each sentence matches the subject given in the hints. Do not  
 486 incorrectly use the subject as the object.

487 *Criteria 6:* Use the correct singular or plural form when referring to the target, which may be a single  
488 object, a pair of objects, or a group of objects.

489 *Criteria 7:* Integrate all relevant information from the hints, noting that some hints may be redundant  
490 or contain errors.

491 Input image: [*Raw Image*].

492 Hint: [*Context-Independent Description*].

### 493 **B.3 Prompt for Rephrasing Referring Expressions**

494 Rewrite the subsequent description while preserving the main information. Utilize varied expressions  
495 and reorganize the sentences if necessary. Begin each sentence with the same subject being referred  
496 to.

497 Description: [*The Referring Expression to be Rephrased*].

### 498 **B.4 Prompt for GPT4-V Evaluation**

499 You are an expert in referring expression comprehension and localization. Your task is to locate the  
500 object in the image based on the provided expression. The coordinates range from the top left (0, 0)  
501 to the bottom right (*Image Width*, *Image Height*). Please provide the bounding box in the format  
502  $(x_0, y_0, x_1, y_1)$ , where  $(x_0, y_0)$  represents the top-left corner and  $(x_1, y_1)$  represents the bottom-right  
503 corner.

504 Expression: [*The Referring Expression*].

## 505 **C More Experiments**

### 506 **C.1 Category-Wise Performance.**

507 Figure 5 presents the per-category performance of the top four models. In Figures 10 and 11, we  
508 show the performance for all 24 models on a per-category basis, with mAcc serving as the metric,  
509 along with the average performance for each model across all categories.

### 510 **C.2 Evaluation on Diverse Data Sources.**

511 Figure 6 illustrates the performance of six models across three subsets, namely “COCO”, “O365-P1”  
512 and “O365-P2”. In Figure 12, the comprehensive results of 24 models across the same three subsets  
513 are displayed.

## 514 **D Limitations and Broad Impacts**

515 Ref-L4 provides a more comprehensive and detailed evaluation of REC capabilities, helping to better  
516 understand and improve the performance of large multimodal models capable of handling the REC  
517 task. The public availability of Ref-L4 and its evaluation code encourages further research and  
518 collaboration, driving innovation and advancements in the field of REC and beyond. While Ref-L4  
519 aims to cover a wide range of scenarios, it may still miss out on specific edge cases or unique contexts  
520 that could be encountered in real-world applications. The detailed and lengthy referring expressions  
521 might pose a challenge for current models, requiring significant advancements in natural language  
522 processing and comprehension capabilities.

## 523 **E Author Statement**

524 The authors of the Ref-L4 benchmark accept full accountability for any rights violations, such as  
525 copyright infringement or other legal breaches. They emphasize that all data included in the Ref-L4



526 dataset adheres to the licensing agreements of the original source datasets. The Ref-L4 benchmark  
527 is made available under the Creative Commons Attribution-NonCommercial 4.0 International (CC  
528 BY-NC 4.0) license. Meticulous attention has been paid to ensure that the dataset upholds the highest  
529 legal and ethical standards. The authors are committed to addressing any issues arising from the use  
530 of this dataset and stand prepared to take necessary actions to resolve them.

## 531 **F Maintenance and Long Term Preservation**

532 To ensure the benchmark remains relevant and useful for evaluating REC models, we will establish  
533 a protocol for regular updates. This includes the addition of new image sets and text annotations  
534 that reflect current trends and challenges in the field. A version control system will be implemented  
535 to track changes and updates to the benchmark. Each version will be documented with detailed  
536 notes on the modifications, including the addition of new data, changes to annotation guidelines, and  
537 improvements based on user feedback. We will utilize reliable cloud storage solutions with multiple  
538 redundancy mechanisms to safeguard against data loss.

## 539 **G Datasheet**

540 The datasheet of our Ref-L4 benchmark can be found in the supplementary material.

## 541 **H Links and Licenses**

542 **Evaluation Code.** The evaluation code is available at [https://github.com/JierunChen/](https://github.com/JierunChen/Ref-L4)  
543 [Ref-L4](https://github.com/JierunChen/Ref-L4).

544 **Benchmark Link.** Ref-L4 is available for download from the Huggingface platform at [https://](https://huggingface.co/datasets/JierunChen/Ref-L4)  
545 [huggingface.co/datasets/JierunChen/Ref-L4](https://huggingface.co/datasets/JierunChen/Ref-L4).

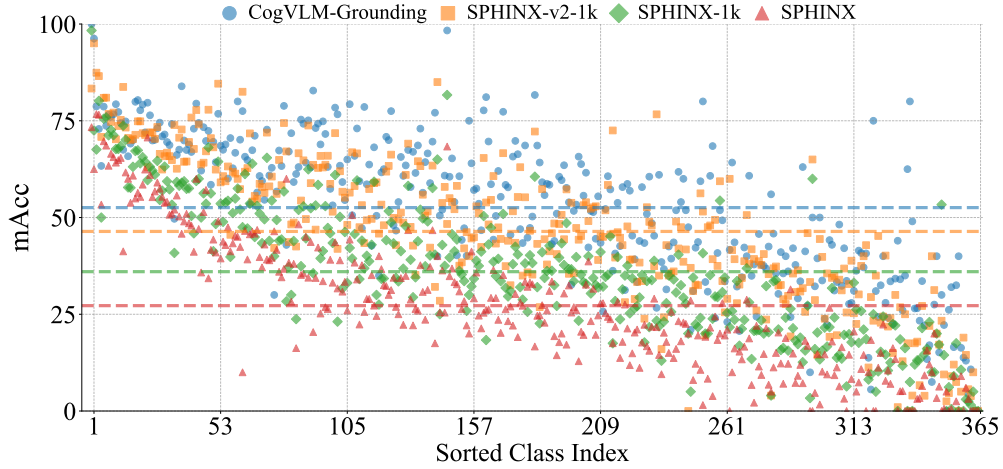
546 **Croissant Metadata.** The croissant format metadata for Ref-L4 can be accessed at [https://](https://huggingface.co/api/datasets/JierunChen/Ref-L4/croissant)  
547 [huggingface.co/api/datasets/JierunChen/Ref-L4/croissant](https://huggingface.co/api/datasets/JierunChen/Ref-L4/croissant).

548 **DOI.** The DOI of Ref-L4 is [10.57967/hf/2388](https://doi.org/10.57967/hf/2388).

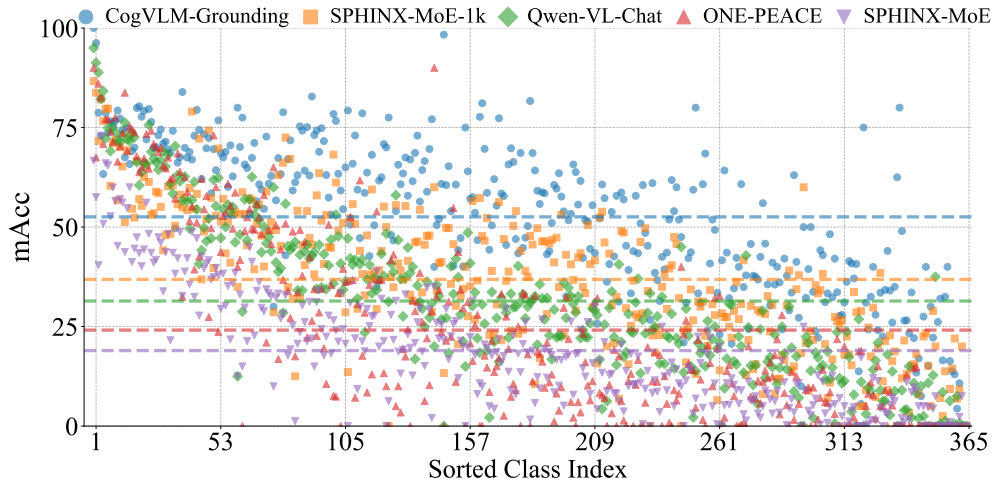
549 **License.** Ref-L4 is distributed under the [Creative Commons Attribution-NonCommercial 4.0 Inter-](https://creativecommons.org/licenses/by-nc/4.0/)  
550 [national \(CC BY-NC 4.0\) license](https://creativecommons.org/licenses/by-nc/4.0/). The images included in Ref-L4 are derived from the following  
551 sources, each governed by their respective licenses:

- 552 • RefCOCO: Licensed under the [Apache-2.0 license](https://www.apache.org/licenses/LICENSE-2.0/).
- 553 • RefCOCO+: Licensed under the [Apache-2.0 license](https://www.apache.org/licenses/LICENSE-2.0/).
- 554 • RefCOCOg: Licensed under the [Creative Commons Attribution 4.0 International \(CC BY](https://creativecommons.org/licenses/by/4.0/)  
555 [4.0\) license](https://creativecommons.org/licenses/by/4.0/).
- 556 • COCO 2014: Licensed under the [Creative Commons Attribution 4.0 International \(CC BY](https://creativecommons.org/licenses/by/4.0/)  
557 [4.0\) license](https://creativecommons.org/licenses/by/4.0/).
- 558 • Objects365: Licensed under the [Creative Commons Attribution 4.0 International \(CC BY](https://creativecommons.org/licenses/by/4.0/)  
559 [4.0\) license](https://creativecommons.org/licenses/by/4.0/).

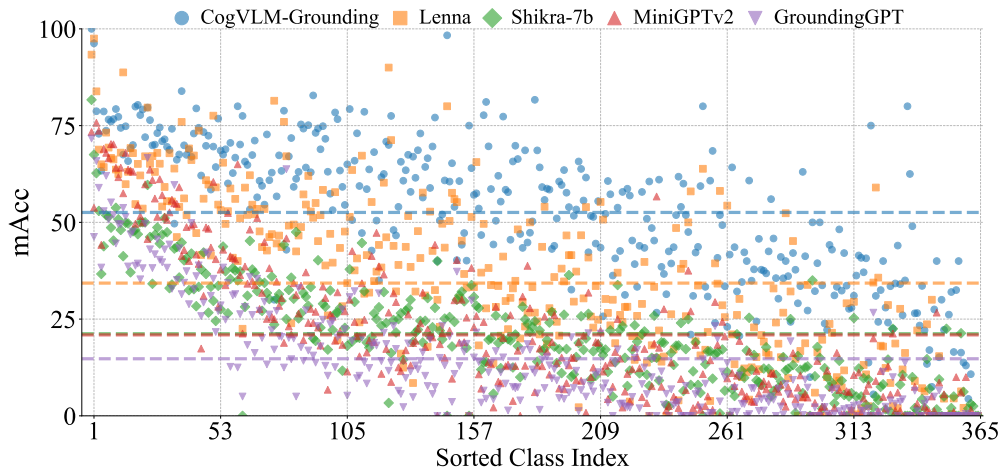




(a) The average performance across all categories (dot lines) for CogVLM-Grounding [62], SPHINX-v2-1k [31], SPHINX-1k [31], and SPHINX1 [31] are 52.56, 46.40, 36.01, and 26.95, respectively.

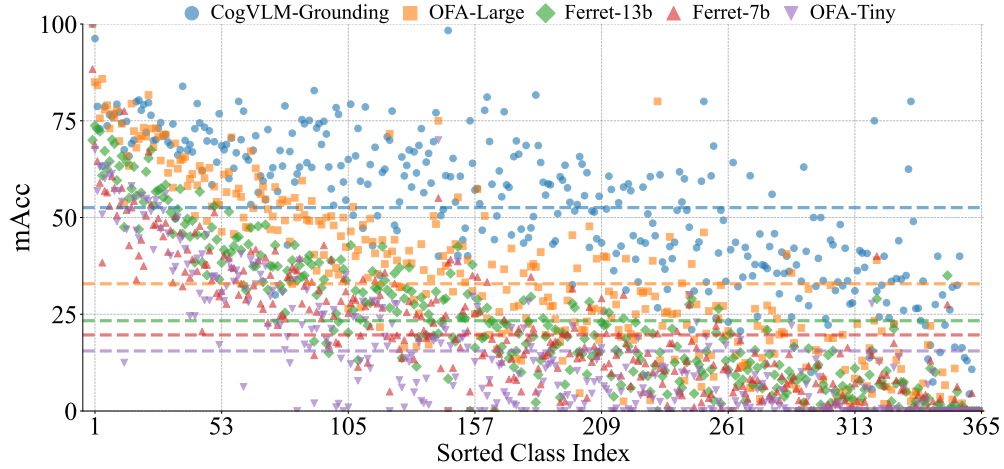


(b) The average performance across all categories (dot lines) for SPHINX-MoE-1k [13], Qwen-VL-Chat [2], ONE-PEACE [60], and SPHINX-MoE [13] are 36.84, 31.41, 24.11, and 18.77, respectively.

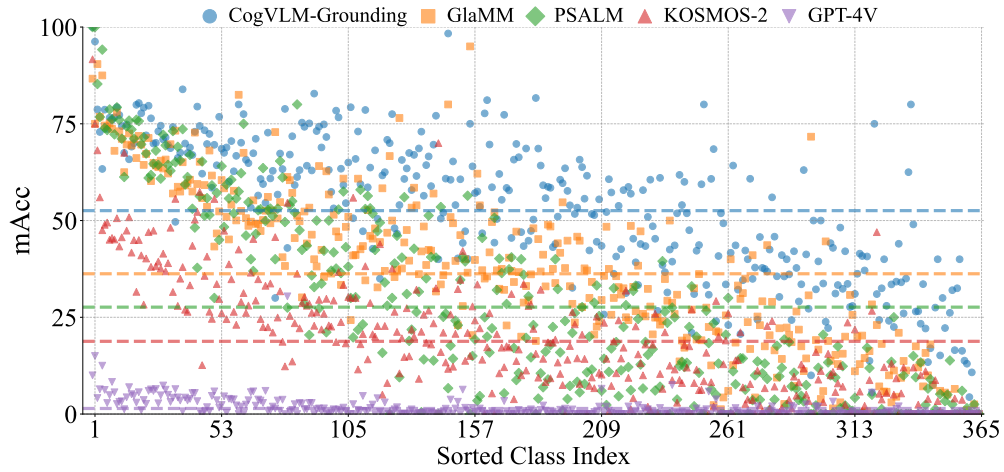


(c) The average performance across all categories (dot lines) for Lenna [66], Shikra-7b [6], MiniGPTv2 [5], and GroundingGPT [29] are 34.30, 21.22, 21.13, and 14.60, respectively.

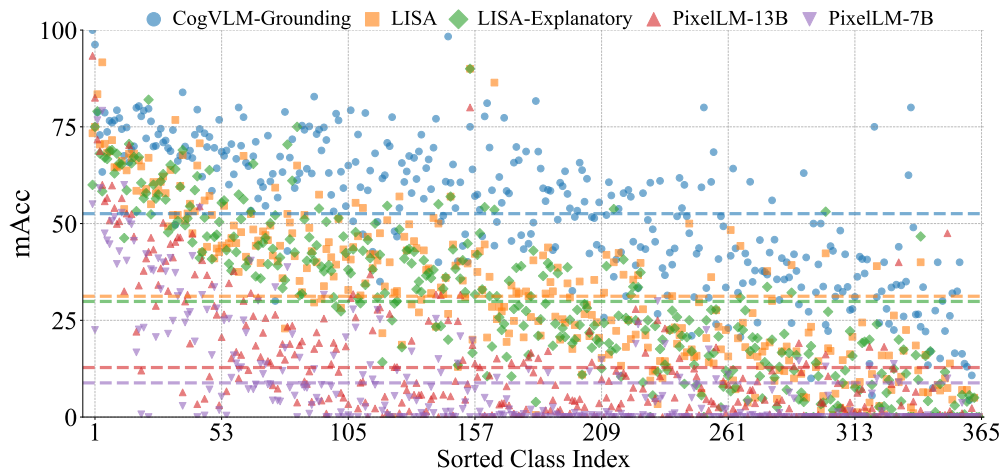
Figure 10: Category-wise performance of 24 models (part-1), sorted in the same order as in Figure 5. We use CogVLM-Grounding as a reference for comparison in each sub-figure.



(a) The average performance across all categories (dot lines) for OFA-Large [59], Ferret-13b [70], Ferret-7b [70] and OFA-Tiny [59] are 32.88, 23.33, 20.27, and 15.37, respectively.



(b) The average performance across all categories (dot lines) for GlaMM [49], PSALM [79], KOSMOS-2 [42] and GPT-4V [39–41] are 36.25, 27.62, 19.37, and 1.42, respectively.



(c) The average performance across all categories (dot lines) for LISA [25], LISA-Explanatory [25], PixelLM-13B [51] and PixelLM-7B [51] are 31.22, 29.87, 13.19, and 8.74, respectively.

Figure 11: Category-wise performance of 24 models (part-2), sorted in the same order as in Figure 5. We use CogVLM-Grounding as a reference for comparison in each sub-figure.

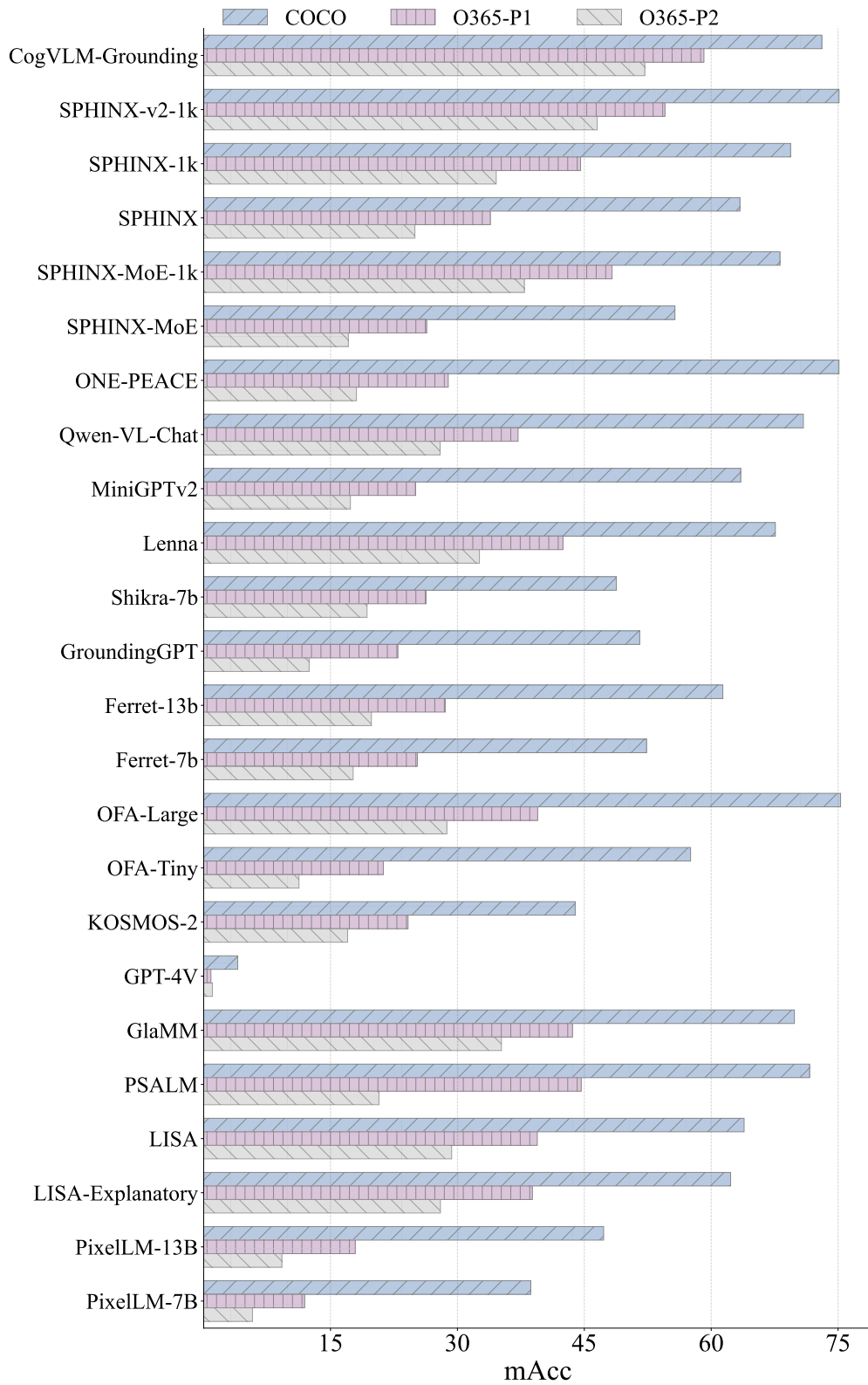


Figure 12: Evaluation of 24 models on various data sources, with mAcc acting as the metric.

560 **Checklist**

- 561 1. For all authors...
- 562 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
563 contributions and scope? [Yes] Refer to Section 1.
- 564 (b) Did you describe the limitations of your work? [Yes] Refer to Section D.
- 565 (c) Did you discuss any potential negative societal impacts of your work? [Yes] Refer to  
566 Section D.
- 567 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
568 them? [Yes]
- 569 2. If you are including theoretical results...
- 570 (a) Did you state the full set of assumptions of all theoretical results? [N/A] No theoretical  
571 results in this work.
- 572 (b) Did you include complete proofs of all theoretical results? [N/A] No theoretical results  
573 in this work.
- 574 3. If you ran experiments (e.g. for benchmarks)...
- 575 (a) Did you include the code, data, and instructions needed to reproduce the main ex-  
576 perimental results (either in the supplemental material or as a URL)? [Yes] Refer to  
577 Sections 3 and 4.
- 578 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
579 were chosen)? [N/A] No training is conducted in this work.
- 580 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
581 ments multiple times)? [N/A] No training is conducted in this work.
- 582 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
583 of GPUs, internal cluster, or cloud provider)? [Yes] NVIDIA A100 (80G) GPUs are  
584 used for evaluation in this work.
- 585 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 586 (a) If your work uses existing assets, did you cite the creators? [Yes] Our benchmark is  
587 derived from RefCOCO, RefCOCO+, RefCOCOg and Objects365.
- 588 (b) Did you mention the license of the assets? [Yes] Refer to Section H.
- 589 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]  
590 We introduce a new benchmark, Ref-L4.
- 591 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
592 using/curating? [N/A]
- 593 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
594 information or offensive content? [N/A]
- 595 5. If you used crowdsourcing or conducted research with human subjects...
- 596 (a) Did you include the full text of instructions given to participants and screenshots, if  
597 applicable? [N/A]
- 598 (b) Did you describe any potential participant risks, with links to Institutional Review  
599 Board (IRB) approvals, if applicable? [N/A]
- 600 (c) Did you include the estimated hourly wage paid to participants and the total amount  
601 spent on participant compensation? [N/A] The benchmark is labeled and reviewed by  
602 the authors of this work.