

---

# Datasheet of the Ref-L4 Benchmark

---

## 1 Motivation

### 1.1 For what purpose was the dataset created?

Existing REC benchmarks such as RefCOCO, RefCOCO+, and RefCOCOg suffer from several limitations such as high labeling error rates, limited vocabulary size and brief referring expressions. The introduced Ref-L4 benchmark is proposed to evaluate modern REC models. Ref-L4 is distinguished by four key features: 1) a substantial sample size with 45,341 annotations, 2) a diverse range of object categories with 365 distinct types and varying instance scales from 30 to 3,767, 3) lengthy referring expressions averaging 24.2 words, and 4) an extensive vocabulary comprising 22,813 unique words. We evaluate a total of 24 models using various evaluation protocols, including accuracy, scale-aware evaluation, and category-wise evaluation.

### 1.2 Who created this dataset and on behalf of which entity?

This dataset was developed by Jierun Chen, Fangyun Wei, Jinjing Zhao, Sizhe Song, Bohuai Wu, Zhuoxuan Peng, S.-H. Gary Chan, and Hongyang Zhang. They conducted this work on behalf of the Hong Kong University of Science and Technology, Microsoft Research Asia, the University of Sydney, and University of Waterloo.

## 2 Composition

### 2.1 How many instances are there in total?

Our final Ref-L4 benchmark encompasses 9,735 images with 45,341 referring expressions, each accurately describing one of the 18,653 unique instances.

### 2.2 What data does each instance consist of?

Each instance is a referring expression that describes a unique target within an image.

### 2.3 Is there a label or target associated with each instance?

Yes, each instance is associated with a bounding box and a detailed referring expression.

### 2.4 Are relationships between individual instances made explicit?

N/A. Each instance is independent.

### 2.5 Are there recommended data splits?

The benchmark is divided into two subsets: a validation set, comprising 30% of the data with 7,231 images, 10,311 instances, and 13,420 referring expressions; and a test set, comprising 70% of the data with 9,467 images, 17,242 instances, and 31,921 referring expressions. Given that our benchmark includes instances from 365 categories, we ensure that each category has at least one sample in both the validation and test sets. While we provide these two splits, we encourage the combined use of both sets for model evaluation.

33 **2.6 Are there any errors, sources of noise, or redundancies in the dataset?**

34 The annotations underwent manual review to reduce errors, especially those caused by hallucina-  
35 tions in descriptions generated by GPT-4V.

36 **2.7 Is the dataset self-contained, or does it link to or rely on external resources?**

37 The dataset is self-contained.

38 **2.8 Does the dataset contain data that might be considered confidential?**

39 No.

40 **2.9 Does the dataset contain data that might be offensive, insulting, threatening, or might  
41 otherwise cause anxiety?**

42 No.

43 **2.10 Does the dataset relate to people?**

44 Our benchmark is constructed using four publicly available datasets: RefCOCO, RefCOCO+, Re-  
45 fCOCOg, and Objects365. These datasets feature images that include the “human” and “person”  
46 categories. However, Ref-L4 is not specifically designed for research focused on human-related  
47 subjects.

48 **2.11 Does the dataset identify any subpopulations?**

49 No.

50 **2.12 Is it possible to identify individuals from the dataset?**

51 No.

52 **2.13 Does the dataset contain data that might be considered sensitive in any way?**

53 No.

54 **3 Collection Process**

55 **3.1 How was the data associated with each instance acquired?**

56 Our benchmark is constructed using four publicly available datasets: RefCOCO, RefCOCO+, Ref-  
57 COCOg, and Objects365.

58 **3.2 If the dataset is a sample from a larger set, what was the sampling strategy?**

59 For the RefCOCO series, we begin by consolidating duplicate images and instances, resulting in  
60 a subset of 6,502 images containing 14,186 unique instances. For Objects365, we select sam-  
61 ples from its testing set based on several criteria: 1) Each image has both height and width  
62 greater than 800 pixels; 2) Each image is sufficiently complex, containing more than 10 cate-  
63 gories and 20 instances; 3) Each instance has a square normalized size  $\sqrt{(hw)/(HW)}$  greater  
64 than 0.05, where  $(h, w)$  represents the instance size and  $(H, W)$  denotes the image size; 4) We  
65 randomly sample  $N$  instances for each of the 365 classes defined in Objects365, with  $N =$   
66  $\min(35, \text{the number of instances for the specific class})$ ; 5) We review and exclude instances with er-  
67 roneous bounding box annotations or those difficult to describe uniquely. For a few rare classes, we  
68 relax criterion-1 to 512 pixels and criterion-2 to 10 objects. Consequently, we collect 3,233 images

69 and 4,467 instances from Objects365. Overall, our Ref-L4 benchmark comprises 9,735 images and  
70 18,653 instances.

### 71 **3.3 Who was involved in the data collection process and how were they compensated?**

72 All authors of this paper were involved in both data collection and the manual review processes.

### 73 **3.4 Over what timeframe was the data collected?**

74 The Ref-L4 benchmark was collected over a period of eight months, spanning from August 2023 to  
75 March 2024.

### 76 **3.5 Were any ethical review processes conducted?**

77 Yes.

### 78 **3.6 Does the dataset relate to people?**

79 N/A.

## 80 **4 Preprocessing/cleaning/labeling**

### 81 **4.1 Was any preprocessing/cleaning/labeling of the data done?**

82 Yes. Given a target instance and its corresponding image, we leverage GPT-4V with human review-  
83 ers in the loop to generate its precise and detailed referring expressions. We manually review all  
84 referring expressions generated by GPT-4V to correct any hallucination issues. We ensure that each  
85 expression uniquely describes the instance and is factual, accurate, and harmless.

### 86 **4.2 Is the software used to preprocess/clean/label the instances available?**

87 Yes, we used GPT-4V, which is publicly available.

## 88 **5 Uses**

### 89 **5.1 Has the dataset been used for any tasks already?**

90 In our study, we used our Ref-L4 benchmark to evaluate 24 REC models.

### 91 **5.2 What (other) tasks could the dataset be used for?**

92 The dataset can be used for evaluating any models, and particularly the large multimodal models  
93 that are capable of handling the REC task.

## 94 **6 Distribution**

### 95 **6.1 Will the dataset be distributed to third parties outside of the entity on behalf of which 96 the dataset was created?**

97 Yes. The Ref-L4 benchmark is available for download from the Huggingface platform at <https://huggingface.co/datasets/JierunChen/Ref-L4>. The DOI is: [10.57967/hf/2388](https://doi.org/10.57967/hf/2388).

### 99 **6.2 When will the dataset be distributed?**

100 The first version of the dataset is scheduled for release in June 2024.

101 **6.3 Will the dataset be distributed under a copyright or other intellectual property (IP)**  
102 **license, and/or under applicable terms of use (ToU)?**

103 The dataset is licensed under [Creative Commons Attribution-NonCommercial 4.0 International \(CC](#)  
104 [BY-NC 4.0\) license](#).

105 **6.4 Do any export controls or other regulatory restrictions apply to the dataset or to**  
106 **individual instances?**

107 No.

## 108 **7 Maintenance**

109 **7.1 Who is supporting/hosting/maintaining the dataset?**

110 The responsibility for maintaining the dataset lies with the authors.

111 **7.2 How can the owner/curator/manager of the dataset be contacted?**

112 The owner of the dataset can be reached at [jchenh@cse.ust.hk](mailto:jchenh@cse.ust.hk) and [fawe@microsoft.com](mailto:fawe@microsoft.com).

113 **7.3 Will the dataset be updated?**

114 The authors of Ref-L4 are dedicated to the ongoing maintenance and preservation of this valuable  
115 dataset. Recognizing its importance for advancing research, we plan to release future updates and ex-  
116 pansion as the dataset is utilized in subsequent studies. Our maintenance strategy includes vigilant  
117 monitoring and prompt resolution of issues identified by the broader research community following  
118 its release. Additionally, we aim to incorporate feedback and contributions from users to ensure the  
119 dataset remains relevant and continues to meet the evolving needs of the academic community.

120 **7.4 Will older versions of the dataset continue to be supported/hosted/maintained?**

121 Yes.

122 **7.5 If others want to extend/augment/build on/contribute to the dataset, is there a**  
123 **mechanism for them to do so?**

124 Yes, the dataset will be released under a Creative Commons (CC) license, allowing others to use,  
125 reproduce, and build upon the data as long as they comply with the terms of the license. This  
126 ensures that researchers and developers can freely contribute to and enhance the dataset, fostering a  
127 collaborative environment for further advancements.

## 128 **8 Reproducibility of the baseline score**

129 We evaluated 24 open-source models. The evaluation code is provided in [https://github.com/](https://github.com/JierunChen/Ref-L4)  
130 [JierunChen/Ref-L4](https://github.com/JierunChen/Ref-L4).

## 131 **9 Reading and using the dataset**

132 Instructions for accessing and using the dataset are available in [https://huggingface.co/](https://huggingface.co/datasets/JierunChen/Ref-L4)  
133 [datasets/JierunChen/Ref-L4](https://huggingface.co/datasets/JierunChen/Ref-L4).

134 **10 Data Format**

135 Each entry in the Ref-L4 dataset includes an image, a bounding box, and a comprehensive referring  
 136 expression annotation. These annotations are stored in parquet format. Table 1 demonstrates the  
 137 key-value structure of the annotation.

Table 1: Key-value structure of the annotation.

Key	Value
id	1
caption	“Within the central picture frame ...”
bbox	$[x, y, w, h]$
bbox_area	10,492.60
bbox_id	“o365_527361”
ori_category_id	“o365_64”
image_id	“o365_922765”
height	741
width	1,024
file_name	“objects365_v2_00922765.jpg”
is_rewrite	true
split	“val”

138 The annotation format primarily consists of the following key-value pairs:

- 139 1. **id**: A unique identifier assigned to each triplet of an image, an instance, and an expression.
- 140 2. **caption**: A comprehensive natural language description of the target instance, offering  
 141 context and detailing specific attributes.
- 142 3. **bbox**: The bounding box coordinates of the target instance, represented as  $[x, y, w, h]$ ,  
 143 where where  $x$  and  $y$  denote the top-left coordinates, and  $w$  and  $h$  signify the width and  
 144 height of the bounding box, respectively.
- 145 4. **bbox\_area** The area of the bounding box.
- 146 5. **bbox\_id**: Unique identifier for the box.
- 147 6. **ori\_category\_id**: Original category identifier.
- 148 7. **image\_id**: Unique identifier for the image.
- 149 8. **height**: Height of the image.
- 150 9. **width**: Width of the image.
- 151 10. **file\_name**: The filename of the image.
- 152 11. **is\_rewrite**: Indicator if the caption is a rewritten version, false for raw caption and true for  
 153 rewritten.
- 154 12. **split**: Benchmark split (“val” or “test”).