

Pseudo-Q: Generating Pseudo Language Queries for Visual Grounding

Haojun Jiang^{1*} Yuanze Lin^{3*†} Dongchen Han¹ Shiji Song¹ Gao Huang^{1,2‡}

¹Tsinghua University, BNRist ²BAAI ³University of Washington

{jhj20, hdc19}@mails.tsinghua.edu.cn, yuanze@uw.edu,

{shijis, gaohuang}@tsinghua.edu.cn

Abstract

Visual grounding, i.e., localizing objects in images according to natural language queries, is an important topic in visual language understanding. The most effective approaches for this task are based on deep learning, which generally require expensive manually labeled image-query or patch-query pairs. To eliminate the heavy dependence on human annotations, we present a novel method, named *Pseudo-Q*, to automatically generate pseudo language queries for supervised training. Our method leverages an off-the-shelf object detector to identify visual objects from unlabeled images, and then language queries for these objects are obtained in an unsupervised fashion with a pseudo-query generation module. Then, we design a task-related query prompt module to specifically tailor generated pseudo language queries for visual grounding tasks. Further, in order to fully capture the contextual relationships between images and language queries, we develop a visual-language model equipped with multi-level cross-modality attention mechanism. Extensive experimental results demonstrate that our method has two notable benefits: (1) it can reduce human annotation costs significantly, e.g., 31% on RefCOCO [65] without degrading original model’s performance under the fully supervised setting, and (2) without bells and whistles, it achieves superior or comparable performance compared to state-of-the-art weakly-supervised visual grounding methods on all the five datasets we have experimented. Code is available at <https://github.com/LeapLabTHU/Pseudo-Q>.

1. Introduction

Visual grounding (VG) task [13, 24, 40, 65] has achieved great progress in recent years, with the advances in both computer vision [16, 20, 21, 25, 26, 46, 56, 57, 59] and natural language processing [4, 14, 41, 50, 53]. It aims to localize the

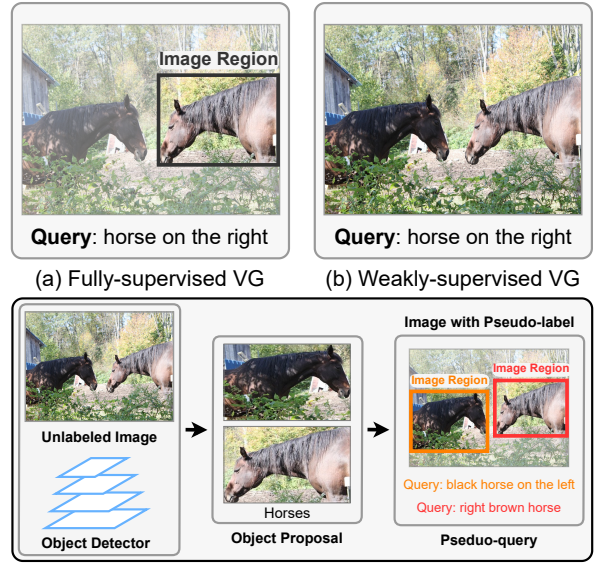


Figure 1. **Comparison with fully and weakly-supervised visual grounding method.** (a) Fully-supervised VG utilizes image region-query pairs as supervision signals. (b) Weakly-supervised VG adopts only language queries. (c) Our **Pseudo-Q** method is free of any task-related annotations.

objects referred by natural language queries, which is essential for various vision-language tasks, e.g., visual question answering [2] and visual commonsense reasoning [67].

Most existing visual grounding methods can be categorized into two types: fully-supervised [8, 13, 22, 23, 33, 35] and weakly-supervised [6, 10, 19, 36, 38, 49, 55, 58]. Although these two lines of works have made remarkable successes, they rely heavily on manually annotated datasets. However, obtaining a large quantity of manual annotations, especially natural language queries, is expensive and time-consuming. To annotate queries, humans need to firstly recognize the visual objects and identify their attributes, and then determine diverse relationships between them on a case-by-case basis, such as spatial (e.g., *left* and *right*), preposition (e.g., *in* and *with*), action (e.g., *throwing some-*

*Equal contribution.

†This work was done during an internship at Tsinghua.

‡Corresponding author.

thing), and comparative (e.g., *smaller* and *bigger*). Among them, *spatial* relation is the most frequently queried one.

To reduce the burden of human annotation, we propose a pseudo language query based approach (**Pseudo-Q**) for visual grounding. Our inspiration comes from previous works [17, 31] that address the high annotation cost issue in image captioning task, by leveraging an unlabelled image set, a sentence corpus, and an off-the-shelf object detector. However, the visual grounding task is more complicated and challenging, as it involves the modelling of relations between objects.

To accurately ground objects by language queries, it’s fundamental to recognize their categories, attributes, and relationships. Thus, when it comes to generating pseudo region-query pairs for an unlabelled image set, we need to focus on three key components: (1) salient **objects (nouns)** which are most likely to be queried, (2) intrinsic **attributes** possessed by the queried objects, and (3) the important **spatial relationships** between the objects. Motivated by [17, 42], we leverage an off-the-shelf object detector [1] to locate the most notable candidates with high confidence, and an attribute classifier [1] to recognize common attributes. However, these detectors are unable to distinguish the spatial relations between objects. Thus, we present a heuristic algorithm to determine the *spatial relationships* between the objects of the same class by comparing their areas and relative coordinates. With these three essential components, pseudo-queries with respect to spatial relations between objects can be generated.

To further improve the performance of our method, we also propose a query prompt module which attentively tailors generated pseudo queries into task-related query templates for visual grounding. For the visual-language model, we put forward a multi-level cross-modality attention mechanism in the fusion module to encourage a deeper fusion between visual and language features.

Extensive experiments have verified the effectiveness of our method. First, in fully supervised manner, it can reduce human annotation costs by 31% without sacrificing original model’s performance on RefCOCO [65]. Second, without bells and whistles, it can obtain superior or comparable performance even compared with state-of-the-art weakly-supervised visual grounding methods on five datasets, including RefCOCO [65], RefCOCO+ [65], RefCOCOg [40], ReferItGame [28] and Flickr30K Entities [44].

In summary, this paper makes three-fold contributions:

- (1) We introduce the first pseudo-query based visual grounding method that deals with the most dominant spatial relationships among objects.
- (2) We propose a query prompt module to specifically tailor pseudo-queries for visual grounding task, and a visual-language model equipped with multi-level cross-modality attention is put forward to fully capture

the contextual relationships of different modalities.

- (3) Extensive experiments demonstrate that our approach can not only dramatically reduce the manual labelling costs without performance sacrifice under the fully supervised condition, but also surpass or achieve comparable performance with state-of-the-art weakly-supervised visual grounding methods.

2. Related Work

2.1. Natural Language Visual Grounding

Visual grounding is a crucial component in vision and language, and it serves as the fundamental of other tasks, such as VQA. Recent visual grounding methods can be summarized into three categories: fully-supervised [8, 13, 22, 23, 33, 35], weakly-supervised [6, 10, 19, 36, 38, 49, 55, 58], and unsupervised [54, 63]. Fully-supervised methods rely heavily on the manual labeled patch-query pairs. Unfortunately, obtaining such sophisticated annotations is expensive and time-consuming. Consequently, weakly-supervised approaches attempt to alleviate the issue by utilizing only image-query pairs. These methods [6, 38] usually leverage a mature object detector to compensate the missing bounding box labels for training. However, annotating the language queries for salient objects in an image is the most laborious part. Thus, unsupervised methods [54, 63] attempt to train a model or directly detect queried objects without any annotations. Our work is also an unsupervised method. However, unlike previous methods, we present a novel method, named Pseudo-Q, to automatically generate pseudo-queries for supervised learning.

2.2. Vision-Language Transformer

Transformer [53] has been firstly proposed to address natural language processing (NLP) tasks. ViT [16] makes the first attempt to apply a transformer for image classification task [12]. Motivated by the success of ViT, DETR [5] and Segmenter [48] further extend the transformer for object detection and segmentation tasks respectively.

There are also many efforts [9, 13, 32, 39, 45, 51], which try to handle visual-language tasks by transformer. TransVG [13] proposes a novel framework with transformer structure for visual grounding task. CLIP [45] and UNITER [9] leverage transformers for jointly learning text and image representations. LXMERT [51] establishes a large-scale transformer to learn cross-modality representation. In this work, we propose a novel multi-level cross-modality attention on the top of the TransVG for cross-modality learning.

2.3. Visual Recognition without Annotation

There have been several works [3, 7, 11, 15, 18, 27, 42, 52, 69] for zero-shot visual tasks. Zero-shot object detection tasks [3, 18] are designed for detecting unseen object

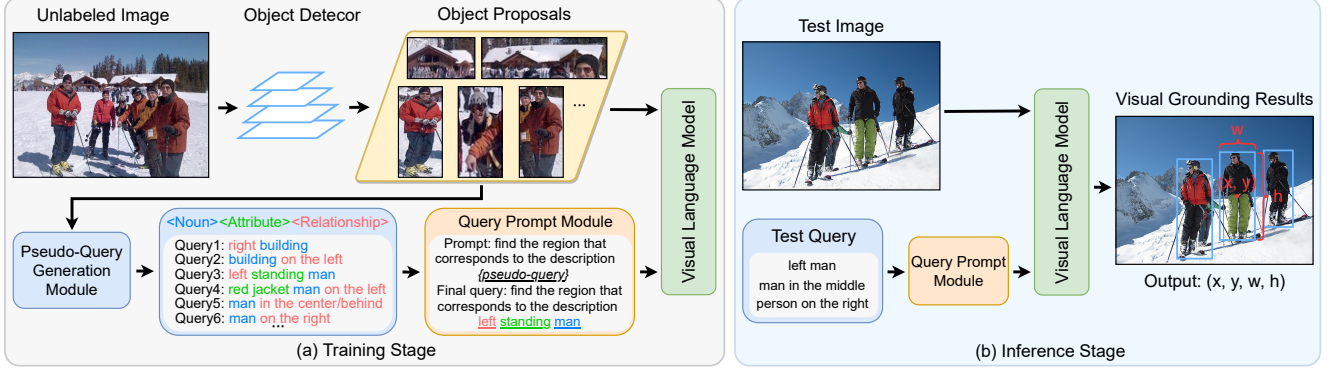


Figure 2. **Overview of our Pseudo-Q method. Better view in color and zoom in.** The proposed approach consists of a pseudo-query generation module, a query prompt module, and a visual-language model. (a) During the training stage, pseudo image region-query pairs are generated to train visual language model. (b) During the inference stage, the test query is filled into the prompt template, and the target object is located by the trained model.

classes whose labels are missing. While zero-shot action recognition task [11, 15, 27] recognizes pre-defined action categories without using action labels. Our work’s emphasis lies in locating object regions without using any task-related annotations, *e.g.*, image regions and queries.

As for zero-shot visual grounding, the pioneering work ZSGNet [47] focuses on query phrases which may contain unseen nouns or object categories. It consists of a language module to encode query features, a visual module to extract image features, and an anchor generator to produce anchors. However, note that the focus of our work is different from ZSGNet, which is proposed for recognizing unseen classes. In addition, ZSGNet utilizes manual annotations while we do not rely on any task-related labels.

3. Method

In this section, we explain our Pseudo-Q method in detail. In Sec. 3.1, we introduce the overview of Pseudo-Q. In Sec. 3.2, we elaborate the pseudo-query generation module. In Sec. 3.3, the details of the task-related query prompt module are shown. Finally, we illustrate the mechanism of our multi-level cross-modality attention in Sec. 3.4.

3.1. Overview

Previous visual grounding methods rely on expensive human annotations, *i.e.*, image region-query pairs for fully-supervised approaches [13, 22, 35] or image-query pairs for weakly-supervised approaches [36, 37, 49]. We firstly propose a pseudo language query based method without using any task-related annotations at training.

Specifically, the Pseudo-Q approach consists of three components, including: (1) pseudo-query generation module, (2) query prompt module, and (3) visual-language model. The illustration of Pseudo-Q is shown in Figure 2. Taking an unlabeled image as an explanation, the detector

can produce several object proposals. Then, these proposals are fed into pseudo-query generation module, which can automatically generate *nouns*, *attributes*, and *relationships* for these proposals. Together with these elements, we can easily create pseudo language queries.

Subsequently, the query prompt module refines created pseudo language queries for visual grounding task. Finally, we propose a visual-language model to fully capture the contextual relationship between the image regions and corresponding pseudo language queries.

3.2. Pseudo-Query Generation

In general, the first step for visual grounding is recognizing the categories of queried objects. However, such a simple grounding strategy leads to ambiguities in complex scenes, *e.g.*, “a talk person on the left” or “a short person on the right”, without understanding their spatial relations or attributes. Thus, to accurately locate visual objects by language queries, a visual grounding model needs to understand queried objects’ categories, attributes, and their relationships. Based on the above analysis, generating pseudo language queries for candidate objects involving three key components: **nouns**, **attributes** and **relationships**.

Nouns. Inspired by works [17, 31, 42], we adopt an off-the-shelf detector [1] to obtain the object proposals. Unlike image classification task where each image contains only one major object, scenes in visual grounding task are more complex due to plenty of candidate objects. While it is natural to select the most salient objects as candidates, such a process requires intensive manual labor which is not available in our setup. Instead, we use detection confidence as a criterion. Concretely, the top- N objects with highest confidence are kept as our proposals. Furthermore, we empirically discover that the detector will focus on a large quantity of tiny objects which are less likely to be queried. Thus, we propose to remove tiny objects before generating proposals.

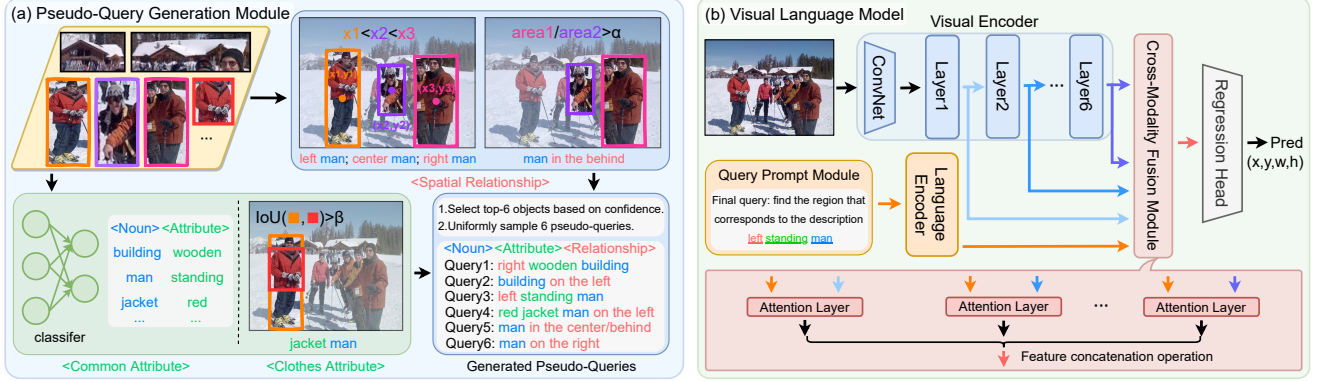


Figure 3. (a) The **pseudo-query generation module** produces spatial relationships and attributes for corresponding objects. (b) The **visual-language model** consists of a visual encoder, a language encoder, and a cross-modality fusion module.

Attributes. They are important semantic cues that help models understand scenes better. We investigate that, in existing datasets [28, 40, 65], common attributes including color, size (*tall*), material (*wooden*) and human state of motion (e.g. *standing* and *walking*), etc. Similar to obtaining the nouns, we take advantage of an off-the-shelf attribute classifier [1] to predict the above common attributes of corresponding objects. In general, one object may have several attributes, such as “a tall person is walking”, and it is ideal to recognize as many attributes as possible. However, limited by the capability of the model, we only keep the attribute with the highest confidence and exceeding a predefined threshold as the final proposal. Furthermore, clothes are also important attributes for a person which can be determined by calculating the IoU value between clothes and person, as shown in Figure 3(a).

Relationships. We observe that *spatial relationship* is one of the most frequently used relations in most existing datasets [40, 65] to distinguish different objects. In order to excavate latent spatial relationships, we propose a heuristic algorithm as shown in Figure 3(a).

In general, spatial relationship can be divided into three dimensions: horizontal (i.e., *left*, *middle*, and *right*), vertical (i.e., *top* and *bottom*) and depth (i.e., *front* and *behind*). Note that each previously generated object proposal is represented by a set of coordinates which naturally embrace spatial information. We can obtain the horizontal and vertical spatial relationships by comparing the center coordinates of objects along with these two dimensions. Meanwhile, to increase the robustness of the algorithm, the numerical difference of two objects’ coordinates in the same dimension is required to be greater than a pre-defined threshold. Finally, we can determine the spatial relations, such as *left*, *right* and *center*, for different visual objects of the same category.

In the depth dimension, we assume that, for the same kind of object, the closer the object is to a camera the larger

the object region. Concretely, we calculate the ratio of the area of the largest object region to the smallest object region and set a threshold to determine whether there is a *front* and *behind* relationship. If the ratio exceeds the threshold, we assign *front* and *behind* relationships to the largest and smallest objects respectively.

Pseudo-queries. After obtaining three key elements, we can generate all possible pseudo-queries for an image following the templates in Appendix. Finally, we sample up to M pseudo image region-query pairs if the number of candidates is greater than M , otherwise, we sample all.

3.3. Query Prompt Module

With the advances of pre-trained language models [4, 14], prompt engineering is proposed to better utilize their learned knowledge at pre-training stage. Inspired by the recent success of prompt engineering in visual-language tasks, e.g., image-language pre-training [45], we propose a query prompt module to excavate the hidden knowledge of pre-trained language model (Sec. 3.4) by refining generated pseudo language queries for visual grounding task.

While the prompt templates proposed in CLIP [45] works well for the image classification task, we empirically find that they are ineffective for the challenging visual grounding task. Consequently, in this work, we explore prompt templates exclusively for visual grounding. Our introduced query prompt module follows certain templates, e.g., “find the region that corresponds to the description {pseudo-query}” or “which region does the text {pseudo-query} describe?”. Such design is specifically tailored for visual grounding task, since the focus of this task lies on locating the regions of referred objects.

3.4. Visual-Language Model

Our visual-language model consists of a visual encoder, a language encoder and a cross-modality fusion module to fuse information from two modalities. The designs of

the visual encoder and the language encoder are following TransVG [13]. We elaborate them in Appendix.

Cross-modality fusion module. Previous method [13] naively utilizes final features of visual and language encoders to acquire cross-modality information. However, such a simple approach is suboptimal, since each level of visual feature possesses valuable semantic information [21, 34]. To be more specific, low-level features usually denote coarse information, *e.g.*, shape and edge, while high-level features can represent finer information, *e.g.*, intrinsic object properties. Thus, we further propose multi-level cross-modality attention (ML-CMA) to thoroughly fuse textual embedding with multi-level visual features.

The mechanism of ML-CMA is shown in Figure 3(b). Features of each visual transformer layer are passed into a cross-modality fusion module with the extracted textual embedding to calculate cross-modality self-attention. Then, we concatenate all updated visual or textual features from different levels respectively, and utilize a fully connected layer to map them into the original dimension. Finally, all features are concatenated and fed into a regression head to predict referred object regions. The regression head composes of three fully connected layers.

4. Experiments

Dataset and setups. Following previous visual grounding methods [13, 60], we evaluate our method on five datasets: RefCOCO [65], RefCOCO+ [65], RefCOCOg [40], ReferItGame [28], and Flickr30K Entities [44]. We follow the same train/val/test splits from [13] for all datasets. The number of training images in these five datasets are 16994, 16992, 24698, 8994, and 29779. Note that we don’t use any manual annotations during the training stage, they are only leveraged for evaluation.

Implementation details. We choose a pre-trained detector [1] and attribute classifier [1] on Visual Genome dataset [30], which contains 1600 object and 400 attribute categories. As we mentioned in Sec. 3.2, we select top- N and sample up to M pseudo-queries for each image. Specifically, on RefCOCO, we select top-3 objects according to the detection confidence and uniformly sample 6 pseudo-queries from all possible candidates. As for RefCOCO+, RefCOCOg, ReferItGame, and Flickr30K Entities, we use top-3 objects/12 pseudo-queries, top-2 objects/4 pseudo-queries, top-6 objects/15 pseudo-queries, and top-7 objects/28 pseudo-queries, respectively.

Training details. All our experiments are conducted under Pytorch framework [43] with 8 RTX3090 GPUs. Our visual-language model is end-to-end optimized with AdamW. The initial learning rate is set to 2.5×10^{-5} for the visual and language encoder and 2.5×10^{-4} for the cross-modality fusion module. The batch size is 256. All the

datasets use cosine learning rate schedule except Flickr30K Entities which adopts exponential decay schedule with 0.85 decay rate. Our model is trained with 10 epochs on RefCOCO, RefCOCOg, and ReferItGame, 20 epochs on RefCOCO+ and Flickr30K Entities. The data augmentations that we utilize are following TransVG [13], *e.g.*, Random-ResizeCrop, RandomHorizontalFlip and ColorJitter.

4.1. Comparison with State-of-the-art Methods

We report comparison results with existing unsupervised [54, 62, 63] and weakly-supervised [38, 49, 55] methods. Note that the weakly-supervised methods are trained with expensive annotated queries. As references, the performance of fully-supervised [13, 60] methods are showed as an upper bound. Specifically, we show the top-1 accuracy (%) results following previous works [38, 55]. The predicted bounding boxes are regarded as correct if the Jaccard overlaps between them and the ground truth are above 0.5.

RefCOCO/RefCOCO+/RefCOCOg. Our method’s performances on RefCOCO, RefCOCO+ and RefCOCOg datasets are reported in Table 1. We compare our method with the existing state-of-the-art unsupervised method CPT [62] and weakly-supervised method DTWREG [49]. Our method can easily surpass CPT by a remarkable margin on all three datasets (*e.g.*, 23.82%/22.15%/23.83% performance improvement on RefCOCO’s val/testA/testB split respectively). When compared with the DTWREG method, our method can achieve better performances on RefCOCO and RefCOCOg datasets. Meanwhile, it can obtain comparable and superior performances on val and testA split of RefCOCO+ dataset. Although we can see that there’s an accuracy gap compared with DTWREG on testB split, our method still gets a large performance gain over CPT. Note that without leveraging any manually labeled queries of RefCOCO+’s training split, our method can still reach considerable performance. All the experiments demonstrate that our generated pseudo-queries can provide effective supervision signals for visual grounding task.

ReferItGame. In Table 2, we show the comparisons with other existing visual grounding methods on ReferItGame dataset. Our method can achieve 43.32% top-1 accuracy, which outperforms all unsupervised and weakly-supervised methods. Especially, compared with the state-of-the-art weakly-supervised method [55], which can achieve 38.39% top-1 accuracy, our method can obtain 4.93% performance improvement without using any annotated labels. These results can demonstrate the superiority of our proposed method.

Flickr30K Entities. The results on Flickr30K Entities dataset is shown in Table 2. It can be observed that our method can still achieve surprising 60.41% top-1 accuracy which outperforms the state-of-the-art weakly-supervised method [38] by 1.14%. Considering the scale of Flickr30K

Table 1. Comparison with state-of-the-art methods on RefCOCO [65], RefCOCO+ [65] and RefCOCOg [40] datasets in terms of top-1 accuracy (%). “Sup.” refers to supervision level: No (without annotation), Weak (only annotated queries), Full (annotated bbox-query pairs). The best two results with supervision level of No and Weak are **bold-faced** and underlined, respectively.

Method	Published on	Sup.	RefCOCO			RefCOCO+			RefCOCOg		
			val	testA	testB	val	testA	testB	val-g	val-u	test-u
CPT [62]	<i>arXiv’21</i>	No	32.20	36.10	30.30	31.90	35.20	28.80	-	<u>36.70</u>	<u>36.50</u>
Ours	<i>CVPR’22</i>		56.02	58.25	54.13	<u>38.88</u>	45.06	32.13	49.82	46.25	47.44
VC [68]	<i>CVPR’18</i>	Weak	-	33.29	30.13	-	34.60	31.58	33.79	-	-
ARN [36]	<i>ICCV’19</i>		34.26	36.43	33.07	34.53	36.01	33.75	33.75	-	-
KPRN [37]	<i>ACMMM’19</i>		35.04	34.74	36.98	35.96	35.24	<u>36.96</u>	33.56	-	-
DTWREG [49]	<i>TPAMI’21</i>		<u>39.21</u>	<u>41.14</u>	<u>37.72</u>	39.18	<u>40.10</u>	38.08	<u>43.24</u>	-	-
MAttNet [64]	<i>CVPR’18</i>	Full	76.65	81.14	69.99	65.33	71.62	56.02	-	66.58	67.27
NMTree [35]	<i>ICCV’19</i>		76.41	81.21	70.09	66.46	72.02	57.52	64.62	65.87	66.44
FAOA [61]	<i>ICCV’19</i>		72.54	74.35	68.50	56.81	60.23	49.60	56.12	61.33	60.36
ReSC [60]	<i>ECCV’20</i>		77.63	80.45	72.30	63.59	68.36	56.81	63.12	67.30	67.20
TransVG [13]	<i>ICCV’21</i>		80.32	82.67	78.12	63.50	68.15	55.63	66.56	67.66	67.44

Table 2. Comparison with state-of-the-art methods on ReferItGame [28] and Flickr30K Entities [44] in terms of top-1 accuracy (%). “Sup.” refers to supervision level: No (without annotation), Weak (only annotated queries), Full (annotated bbox-query pairs). The best two results with supervision level of No and Weak are **bold-faced** and underlined, respectively.

Method	Published on	Sup.	ReferIt	Flickr30K
Yeh <i>et al.</i> [63]	<i>CVPR’18</i>	No	36.93	20.91
Wang <i>et al.</i> [54]	<i>ICCV’19</i>		26.48	50.49
Ours	<i>CVPR’22</i>		43.32	60.41
Chen <i>et al.</i> [6]	<i>CVPR’18</i>	Weak	33.67	46.61
Zhao <i>et al.</i> [70]	<i>CVPR’18</i>		33.10	13.61
Liu <i>et al.</i> [36]	<i>ICCV’19</i>		26.19	-
Gupta <i>et al.</i> [19]	<i>ECCV’20</i>		-	51.67
Liu <i>et al.</i> [38]	<i>CVPR’21</i>		37.68	<u>59.27</u>
Wang <i>et al.</i> [55]	<i>CVPR’21</i>		<u>38.39</u>	53.10
Kovvuri <i>et al.</i> [29]	<i>ACCV’18</i>	Full	59.13	72.83
Yu <i>et al.</i> [66]	<i>IJCAI’18</i>		63.00	73.30
Yang <i>et al.</i> [61]	<i>ICCV’19</i>		60.67	68.71
Yang <i>et al.</i> [60]	<i>ECCV’20</i>		64.60	69.28
Deng <i>et al.</i> [13]	<i>ICCV’21</i>		69.76	78.47

Entities which consists of 427K manually annotated referred expressions, our method still achieves remarkable performance without any training label. As for other methods without using manual labels, our method can easily surpass [63] and [54] with 39.50% and 9.92% absolute performance improvement, respectively.

Explanations of the gain over weakly-supervised methods. First, the core of visual grounding task is learning the *correspondence* between visual and linguistic modalities which relies heavily on the correct mapping between image regions and queries inside training data. A *key difference* between our approach and weakly-supervised methods is that we can generate corresponding queries for the detected object which guarantees the *correctness of map-*

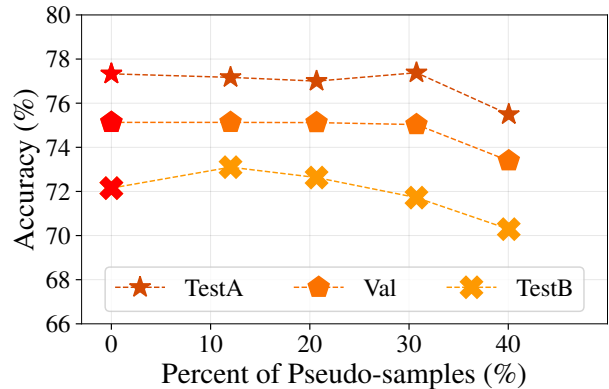


Figure 4. Experiments of reducing the manual labeling cost on RefCOCO [65]. We replace the manual labels whose queries contain spatial relationships with our pseudo-samples.

ping between two modalities. Although weakly-supervised methods have annotated queries, they lack key supervision signals that are the region-level correspondence between two modalities. Second, our model jointly optimizes the features from two modalities which allows the model to learn a better correspondence while weakly-supervised methods [38, 49, 55] only update the language model leaving the visual model fixed.

4.2. Improving the Efficiency of Manual Labeling.

In Figure 4, we perform experiments with the same hyper-parameters as Sec. 4 on RefCOCO [65] to verify the effectiveness of our pseudo-samples, *i.e.*, pseudo image region-query pairs, by replacing the manually annotated labels whose queries contain spatial relationships. The baseline is our model trained in a fully-supervised manner. Note that the query prompt module is not applied in this experiment. As we can see, compared to the fully-supervised set-

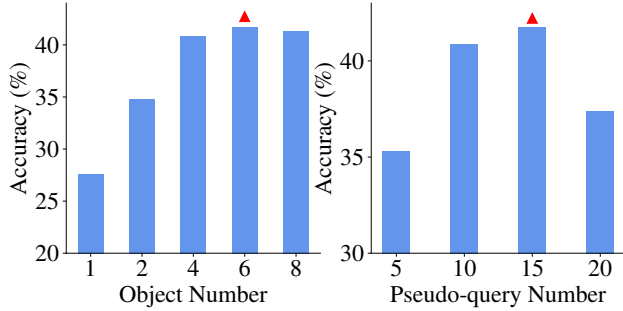


Figure 5. Left: Ablation of object number. Right: Ablation of pseudo-query number. Both are conducted on ReferItGame [40].

ting, substituting 12.01%, 20.68%, and 30.75% manually annotated labels with our generated pseudo region-query pairs do not degrade the original performance. In such a situation, about **31%** of human annotation costs can be reduced. Consequently, our method can be utilized to automatically annotate one of the dominant components, *i.e.* spatial relationship, in language queries, which significantly improves the efficiency of manual labeling.

4.3. Ablation Study

In this section, we conduct extensive ablation experiments to demonstrate the effectiveness of each proposed component and the rationality of hyper-parameters setting. Most of the following experiments are conducted on ReferItGame [28] dataset and we report the top-1 accuracy. The model is trained with the same hyper-parameters as Sec. 4.

Number of nouns. We investigate the impact of utilizing different number of nouns (objects) in Figure 5(a). Increasing the number of nouns can produce more pseudo samples which boosts the performance of our model, as shown in Figure 5(a). In our experiments, we use the detection confidence as a criterion to select salient objects. If the number of nouns is too large, the likelihood of detecting low confidence objects which are inconspicuous will grow. We empirically find that, on ReferItGame dataset, the model reaches its peak performance when the number of nouns is 6. Once the number of nouns exceeds 6, the performance starts to degrade. Thus, we use the top-6 object proposals on the ReferItGame dataset.

Number of pseudo-queries. Another essential factor is the number of pseudo-queries. We study the influence of sampling different number of pseudo-queries in Figure 5(b). The candidates of pseudo-queries are generated following templates in Appendix. As shown in Figure 5(b), our model performs best when the sampling number of pseudo-queries is 15. If the sampling number is too small, we will miss plenty of useful candidates which hinders the improvement of model performance. Meanwhile, note that not every candidate provides the correct supervision signal. Thus, overly

Table 3. Ablations of each component on RefCOCO [65] and ReferItGame [40]. “Attr” and “Rela” denote attribute and relationship, respectively. “Prompt” represents the query prompt module. “ML-CMA” means the proposed multi-level cross-modality attention.

Noun	Attr	Rela	ML-CMA	Prompt	RefCOCO	ReferIt
✓					22.04	27.91
✓	✓				31.30 (↑9.26)	31.33 (↑3.42)
✓		✓			48.71 (↑26.67)	39.26 (↑11.35)
✓	✓	✓			53.39 (↑31.35)	40.37 (↑12.46)
✓	✓	✓	✓		55.16 (↑1.77)	41.72 (↑1.35)
✓	✓	✓	✓	✓	56.02 (↑2.63)	43.32 (↑2.95)

sampling candidates will also hurt the performance. Finally, we sample up to 15 pseudo-queries.

Effectiveness of integrating attributes. We empirically support the effectiveness of introducing attributes into pseudo-queries by comparing them with those lacking attributes. As shown in Table 3, generating pseudo-queries with nouns and attributes clearly surpasses the one that only has nouns on RefCOCO and ReferItGame. Moreover, adding attributes into pseudo-queries with nouns and relationships can further boost the performance. Thus, we demonstrate that incorporating the attribute into pseudo-queries helps models to comprehend the scenes better.

Effectiveness of generating relationships. As we mentioned in Sec. 3.2, spatial relationship is the most essential component. With only nouns, models are still far from comprehensively understanding scenes. The ablation study of generating relationships that supports our proposition is reported in Table 3. Compared with pseudo-queries with only nouns, generating relationships with our method outperforms it overwhelmingly by **26.67%** on RefCOCO. In sum, experimental results show that incorporating spatial relationships into pseudo-queries can significantly enhance the model’s capability of understanding scenes.

Effectiveness of query prompts. In Table 3, we show that prompts help to excavate the hidden knowledge of the pre-trained language model, and consequently, boost the performance. On ReferItGame, the well-designed prompt “*which region does the text {pseudo-query} describe?*” improves the performance by 1.60%. Meanwhile, on RefCOCO, the prompt “*find the region that corresponds to the description {pseudo-query}*” improves the performance by 0.86%. On the other hand, we find that hand-designed prompts are not robust enough across all the datasets.

Effectiveness of cross-modality fusion module. We further investigate the contribution of the cross-modality fusion module in Table 3. On the basis of the pseudo-query generation module, the proposed cross-modality fusion module can further improve the performance by 1.77% and 1.35% on RefCOCO and ReferItGame, respectively.



Figure 6. Four visualization examples of detection results. The red bounding boxes and queries are ground truth. The green bounding boxes are detected by the model that trained on pseudo-samples generated with nouns, attributes, and relationships. The blue bounding boxes are detected by the model that trained on pseudo-samples generated with only nouns.



Figure 7. Four visualization examples of pseudo-sample generated by our method. We use blue and green to distinguish two objects.

4.4. Qualitative Analysis

To further figure out the importance of spatial relationships and attributes, in Figure 6, we show four detection examples of our models trained on generated pseudo-queries with or without spatial relationships and attributes on RefCOCO dataset. In the first two examples, we can easily observe that the model trained with relationships locates target objects much better than the one trained without relationship component. In the last two example, the key factor to locate the queried man is leveraging the attributes “blue” and “standing”. Obviously, with the above analysis, we can conclude that the proposed spatial relationship and attribute play an essential role in accurately grounding referred objects with given language queries. In addition, we also display four generated pseudo region-query pairs on RefCOCO dataset in Figure 7.

4.5. Limitation

Although our method achieves superior performances on five datasets, there are still two limitations. First, when it comes to pseudo language queries, there may be some incorrect queries which harms final performance. Second, the generated pseudo-queries are simple, other relationships can be explored in the future.

5. Conclusion

In this paper, we make the first attempt to introduce a pseudo-query based visual grounding method called Pseudo-Q. Firstly, we propose a pseudo-query generation module to automatically produce pseudo region-query pairs for supervised training. Then, we present a query prompt module, so that generated pseudo language queries can be tailored specifically for visual grounding task. Finally, in order to sufficiently model the relationships between visual regions and language-queries, we develop a visual-language model equipped with multi-level cross-modality attention. Extensive experiments have shown that our method can not only achieve superior performances on five datasets, but also dramatically reduce manual labeling costs.

Acknowledgement

This work is supported in part by the National Science and Technology Major Project of the Ministry of Science and Technology of China under Grants 2018AAA0100701, the National Natural Science Foundation of China under Grants 61906106 and 62022048, the Guoqiang Institute of Tsinghua University and Beijing Academy of Artificial Intelligence.

References

- [1] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. Bottom-up and top-down attention for image captioning and visual question answering. In *CVPR*, 2018. 2, 3, 4, 5
- [2] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *ICCV*, 2015. 1
- [3] Ankan Bansal, Karan Sikka, Gaurav Sharma, Rama Chellappa, and Ajay Divakaran. Zero-shot object detection. In *ECCV*, 2018. 2
- [4] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. In *NeurIPS*, 2020. 1, 4
- [5] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020. 2, 11
- [6] Kan Chen, Jiyang Gao, and Ram Nevatia. Knowledge aided consistency for weakly supervised phrase grounding. In *CVPR*, 2018. 1, 2, 6
- [7] Shizhe Chen and Dong Huang. Elaborative rehearsal for zero-shot action recognition. In *ICCV*, 2021. 2
- [8] Xinpeng Chen, Lin Ma, Jingyuan Chen, Zequn Jie, Wei Liu, and Jiebo Luo. Real-time referring expression comprehension by single-stage grounding network. *arXiv preprint arXiv:1812.03426*, 2018. 1, 2
- [9] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Universal image-text representation learning. In *ECCV*, 2020. 2
- [10] Samyak Datta, Karan Sikka, Anirban Roy, Karuna Ahuja, Devi Parikh, and Ajay Divakaran. Align2ground: Weakly supervised phrase grounding guided by image-caption alignment. In *ICCV*, 2019. 1, 2
- [11] Berkan Demirel, Ramazan Gokberk Cinbis, and Nazli Ikizler-Cinbis. Attributes2classname: A discriminative model for attribute-based unsupervised zero-shot learning. In *ICCV*, 2017. 2, 3
- [12] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 2, 11
- [13] Jiajun Deng, Zhengyuan Yang, Tianlang Chen, Wengang Zhou, and Houqiang Li. Transvg: End-to-end visual grounding with transformers. In *ICCV*, 2021. 1, 2, 3, 5, 6
- [14] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 1, 4, 11
- [15] Mandar Dixit, Yunsheng Li, and Nuno Vasconcelos. Semantic fisher scores for task transfer: Using objects to classify scenes. *IEEE TPAMI*, 2019. 2, 3
- [16] 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. 1, 2
- [17] Yang Feng, Lin Ma, Wei Liu, and Jiebo Luo. Unsupervised image captioning. In *CVPR*, 2019. 2, 3
- [18] Xiuye Gu, Tsung-Yi Lin, Weicheng Kuo, and Yin Cui. Zero-shot detection via vision and language knowledge distillation. *arXiv preprint arXiv:2104.13921*, 2021. 2
- [19] Tanmay Gupta, Arash Vahdat, Gal Chechik, Xiaodong Yang, Jan Kautz, and Derek Hoiem. Contrastive learning for weakly supervised phrase grounding. In *ECCV*, 2020. 1, 2, 6
- [20] Yizeng Han, Gao Huang, Shiji Song, Le Yang, Yitian Zhang, and Haojun Jiang. Spatially adaptive feature refinement for efficient inference. *IEEE TIP*, 2021. 1
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 1, 5, 11
- [22] Richang Hong, Daqing Liu, Xiaoyu Mo, Xiangnan He, and Hanwang Zhang. Learning to compose and reason with language tree structures for visual grounding. *IEEE TPAMI*, 2019. 1, 2, 3
- [23] Ronghang Hu, Marcus Rohrbach, Jacob Andreas, Trevor Darrell, and Kate Saenko. Modeling relationships in referential expressions with compositional modular networks. In *CVPR*, 2017. 1, 2
- [24] Ronghang Hu, Huazhe Xu, Marcus Rohrbach, Jiashi Feng, Kate Saenko, and Trevor Darrell. Natural language object retrieval. In *CVPR*, 2016. 1
- [25] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *CVPR*, 2017. 1
- [26] Gao Huang, Yulin Wang, Kangchen Lv, Haojun Jiang, Wenhui Huang, Pengfei Qi, and Shiji Song. Glance and focus networks for dynamic visual recognition. *arXiv preprint arXiv:2201.03014*, 2022. 1
- [27] Mihir Jain, Jan C Van Gemert, Thomas Mensink, and Cees GM Snoek. Objects2action: Classifying and localizing actions without any video example. In *ICCV*, 2015. 2, 3
- [28] Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. Referitgame: Referring to objects in photographs of natural scenes. In *EMNLP*, 2014. 2, 4, 5, 6, 7
- [29] Rama Kovvuri and Ram Nevatia. Pirc net: Using proposal indexing, relationships and context for phrase grounding. In *ACCV*, 2018. 6
- [30] 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. *IJCV*, 2017. 5
- [31] Iro Laina, Christian Rupprecht, and Nassir Navab. Towards unsupervised image captioning with shared multimodal embeddings. In *ICCV*, 2019. 2, 3

- [32] Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, et al. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *ECCV*, 2020. 2
- [33] Yue Liao, Si Liu, Guanbin Li, Fei Wang, Yanjie Chen, Chen Qian, and Bo Li. A real-time cross-modality correlation filtering method for referring expression comprehension. In *CVPR*, 2020. 1, 2
- [34] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *CVPR*, 2017. 5
- [35] Daqing Liu, Hanwang Zhang, Feng Wu, and Zheng-Jun Zha. Learning to assemble neural module tree networks for visual grounding. In *ICCV*, 2019. 1, 2, 3, 6
- [36] Xuejing Liu, Liang Li, Shuhui Wang, Zheng-Jun Zha, Dechao Meng, and Qingming Huang. Adaptive reconstruction network for weakly supervised referring expression grounding. In *ICCV*, 2019. 1, 2, 3, 6
- [37] Xuejing Liu, Liang Li, Shuhui Wang, Zheng-Jun Zha, Li Su, and Qingming Huang. Knowledge-guided pairwise reconstruction network for weakly supervised referring expression grounding. In *ACMMM*, 2019. 3, 6
- [38] Yongfei Liu, Bo Wan, Lin Ma, and Xuming He. Relation-aware instance refinement for weakly supervised visual grounding. In *CVPR*, 2021. 1, 2, 5, 6
- [39] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *NeurIPS*, 2019. 2
- [40] Junhua Mao, Jonathan Huang, Alexander Toshev, Oana Camburu, Alan L Yuille, and Kevin Murphy. Generation and comprehension of unambiguous object descriptions. In *CVPR*, 2016. 1, 2, 4, 5, 6, 7
- [41] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *NeurIPS*, 2013. 1
- [42] Jinwoo Nam, Daechul Ahn, Dongyeop Kang, Seong Jong Ha, and Jonghyun Choi. Zero-shot natural language video localization. In *ICCV*, 2021. 2, 3
- [43] 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. In *NeurIPS*, 2018. 5
- [44] Bryan A Plummer, Liwei Wang, Chris M Cervantes, Juan C Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In *ICCV*, 2015. 2, 5, 6
- [45] 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 *ICML*, 2021. 2, 4
- [46] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NeurIPS*, 2015. 1
- [47] Arka Sadhu, Kan Chen, and Ram Nevatia. Zero-shot grounding of objects from natural language queries. In *ICCV*, 2019. 3
- [48] Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. Segmenter: Transformer for semantic segmentation. In *ICCV*, 2021. 2
- [49] Mingjie Sun, Jimin Xiao, Enggee Lim, Si Liu, and John Yannis Goulermas. Discriminative triad matching and reconstruction for weakly referring expression grounding. *IEEE TPAMI*, 2021. 1, 2, 3, 5, 6
- [50] Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks. In *NeurIPS*, 2014. 1
- [51] Hao Tan and Mohit Bansal. Lxmert: Learning cross-modality encoder representations from transformers. In *EMNLP*, 2019. 2
- [52] Damien Teney and Anton van den Hengel. Zero-shot visual question answering. *arXiv preprint arXiv:1611.05546*, 2016. 2
- [53] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 1, 2
- [54] Josiah Wang and Lucia Specia. Phrase localization without paired training examples. In *ICCV*, 2019. 2, 5, 6
- [55] Liwei Wang, Jing Huang, Yin Li, Kun Xu, Zhengyuan Yang, and Dong Yu. Improving weakly supervised visual grounding by contrastive knowledge distillation. In *CVPR*, 2021. 1, 2, 5, 6
- [56] Yulin Wang, Zhaoxi Chen, Haojun Jiang, Shiji Song, Yizeng Han, and Gao Huang. Adaptive focus for efficient video recognition. In *ICCV*, 2021. 1
- [57] Yulin Wang, Yang Yue, Yuanze Lin, Haojun Jiang, Zihang Lai, Victor Kulikov, Nikita Orlov, Humphrey Shi, and Gao Huang. Adafocus v2: End-to-end training of spatial dynamic networks for video recognition. In *CVPR*, 2022. 1
- [58] Fanyi Xiao, Leonid Sigal, and Yong Jae Lee. Weakly-supervised visual grounding of phrases with linguistic structures. In *CVPR*, 2017. 1, 2
- [59] Le Yang, Haojun Jiang, Ruojin Cai, Yulin Wang, Shiji Song, Gao Huang, and Qi Tian. Condensenet v2: Sparse feature reactivation for deep networks. In *CVPR*, 2021. 1
- [60] Zhengyuan Yang, Tianlang Chen, Liwei Wang, and Jiebo Luo. Improving one-stage visual grounding by recursive sub-query construction. In *ECCV*, 2020. 5, 6
- [61] Zhengyuan Yang, Boqing Gong, Liwei Wang, Wenbing Huang, Dong Yu, and Jiebo Luo. A fast and accurate one-stage approach to visual grounding. In *ICCV*, 2019. 6
- [62] Yuan Yao, Ao Zhang, Zhengyan Zhang, Zhiyuan Liu, Tat-Seng Chua, and Maosong Sun. Cpt: Colorful prompt tuning for pre-trained vision-language models. *arXiv preprint arXiv:2109.11797*, 2021. 5, 6
- [63] Raymond A Yeh, Minh N Do, and Alexander G Schwing. Unsupervised textual grounding: Linking words to image concepts. In *CVPR*, 2018. 2, 5, 6
- [64] Licheng Yu, Zhe Lin, Xiaohui Shen, Jimei Yang, Xin Lu, Mohit Bansal, and Tamara L Berg. Mattnet: Modular at-

tention network for referring expression comprehension. In *CVPR*, 2018. 6

- [65] Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. Modeling context in referring expressions. In *ECCV*, 2016. 1, 2, 4, 5, 6, 7, 11
- [66] Zhou Yu, Jun Yu, Chenchao Xiang, Zhou Zhao, Qi Tian, and Dacheng Tao. Rethinking diversified and discriminative proposal generation for visual grounding. In *IJCAI*, 2018. 6
- [67] Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. From recognition to cognition: Visual commonsense reasoning. In *CVPR*, 2019. 1
- [68] Hanwang Zhang, Yulei Niu, and Shih-Fu Chang. Grounding referring expressions in images by variational context. In *CVPR*, 2018. 6
- [69] Lingling Zhang, Xiaojun Chang, Jun Liu, Minnan Luo, Sen Wang, Zongyuan Ge, and Alexander Hauptmann. Zstad: Zero-shot temporal activity detection. In *CVPR*, 2020. 2
- [70] Fang Zhao, Jianshu Li, Jian Zhao, and Jiashi Feng. Weakly supervised phrase localization with multi-scale anchored transformer network. In *CVPR*, 2018. 6

Appendix

A. Statistics of RefCOCO Dataset

In Figure 8, we show the statistics of the training set of RefCOCO [65] dataset to demonstrate spatial relationship is one of the dominant components in language queries. As we can see, spatial relationships exists in almost 60% of queries. Furthermore, the most common spatial relationships in RefCOCO are *left* and *right*. In addition, other spatial relationships, *i.e.*, *middle*, *front*, *top*, and *bottom*, are also frequently found in language queries.

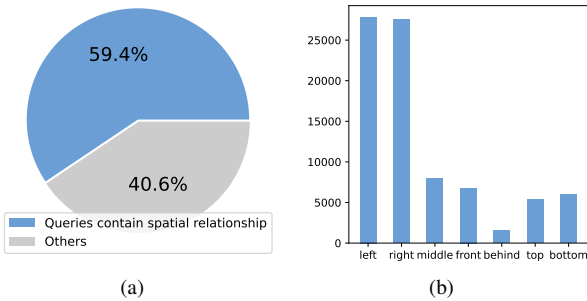


Figure 8. **Statistics of the training set of RefCOCO [65] dataset.** (a): The percent of language queries that contain spatial relationships. (b): The number of different spatial relationships.

B. Pseudo-Query Templates

Our pseud-queries are generated following the templates shown in Table. 4. All the possible templates is considered in our method for the purpose of obtaining as many candidate pseudo-samples as possible. Honestly, this strategy will inevitably produce some ungrammatical pseudo-samples. Our approach is similar to all the pseudo-label

based methods, such as semi-supervised learning, which can’t guarantee every single pseudo-query is correct. Overall, these pseudo-queries provide valuable supervision signals and eventually benefit the training of the model.

Table 4. Pseudo-query templates. *Attr* and *Rela* represents attribute and relationship, respectively.

Pseudo Query Template	Example
{ <i>Noun</i> }	“man”, “building” etc.
{ <i>Noun</i> } { <i>Attr</i> } { <i>Attr</i> } { <i>Noun</i> }	“man standing” etc. “talk man”, “wooden building” etc.
{ <i>Noun</i> } { <i>Rela</i> } { <i>Rela</i> } { <i>Noun</i> }	“man on the right” etc. “center man”, “left building” etc.
{ <i>Noun</i> } { <i>Attr</i> } { <i>Rela</i> } { <i>Noun</i> } { <i>Rela</i> } { <i>Attr</i> }	“man standing on the right” etc. “man right standing” etc.
{ <i>Attr</i> } { <i>Noun</i> } { <i>Rela</i> } { <i>Attr</i> } { <i>Rela</i> } { <i>Noun</i> }	“standing man on the right” etc. “standing right man” etc.
{ <i>Rela</i> } { <i>Noun</i> } { <i>Attr</i> } { <i>Rela</i> } { <i>Attr</i> } { <i>Noun</i> }	“right man standing” etc. “right standing man” etc.

C. Visual-Language Model

In this section, we provide more details about the architecture of the visual encoder and the language encoder.

In the visual encoder, a CNN backbone and a transformer-based network are stacked sequentially for image feature extraction. The CNN backbone is a ResNet-50 model [21] pre-trained on ImageNet [12], and the transformer-based network is the encoder part of DETR network [5] which consists of six transformer layers. Moreover, the pre-trained weights of DETR are utilized for initialization. The output feature maps of the ResNet-50 are fed into a 1×1 convolutional layer for dimension reduction. Then, they are flattened into 1D vectors for the transformer network.

In the language encoder, a token embedding layer and a linguistic transformer are employed to extract textual features. Specifically, the token embedding layer is leveraged to convert the discrete words into continuous language vectors. Since BERT [14] has been successfully applied for text feature extraction, the BERT architecture which has 12 transformer layers is adopted as the linguistic transformer.