# Iterative Robust Visual Grounding with Masked Reference based Centerpoint Supervision

Menghao Li[1] *    Chunlei Wang[1] *    Wenquan Feng[1]    Shuchang Lyu[1]

Guangliang Cheng[2]✉    Xiangtai Li[3]    Binghao Liu[1]    Qi Zhao[1]✉

[1] Beihang University [2] University of Liverpool [3] S-Lab, Nanyang Technological University

{sy2102227, wcl_buaa, buaafwq, lyushuchang, liubinghao, zhaoqi}@buaa.edu.cn

{Guangliang.Cheng}@liverpool.ac.uk    {xiangtai.li}@ntu.edu.sg

## Abstract

*Visual Grounding (VG) aims at localizing target objects from an image based on given expressions and has made significant progress with the development of detection and vision transformer. However, existing VG methods tend to generate **false-alarm** objects when presented with inaccurate or irrelevant descriptions, which commonly occur in practical applications. Moreover, existing methods fail to capture fine-grained features, accurate localization, and sufficient context comprehension from the whole image and textual descriptions. To address both issues, we propose an Iterative Robust Visual Grounding (IR-VG) framework with Masked Reference based Centerpoint Supervision (MRCS). The framework introduces iterative multi-level vision-language fusion (IMVF) for better alignment. We use MRCS to ahieve more accurate localization with point-wised feature supervision. Then, to improve the robustness of VG, we also present a multi-stage false-alarm sensitive decoder (MFSD) to prevent the generation of false-alarm objects when presented with inaccurate expressions. Extensive experiments demonstrate that IR-VG achieves new state-of-the-art (SOTA) results, with improvements of 25% and 10% compared to existing SOTA approaches on the two newly proposed robust VG datasets. Moreover, the proposed framework is also verified effective on five **regular** VG datasets. Codes and models will be publicly at https://github.com/cv516Buaa/IR-VG.*

## 1. Introduction

Visual Grounding (VG) is a crucial computer vision task gaining significant attention due to its potential for enabling practical applications such as robot navigation [11] and vi-

---

*Contribute Equally.



a black and brown cow walking through the ocean water (a)

A man wearing nike shoes in a bright neon blue top is playing tennis

A cat standing on the wash basin with its head up (b)

The front edge of a tan scooter with a carrying container on it (c)

Figure 1. Weaknesses illustration of the existing VG approaches. Green and blue boxes represent groundtruths and prediction. Green boxes and blue boxes represent the ground truths and predictions, respectively. (a) Failure cases occur when an irrelevant or inaccurate description is provided. (b) Fine-grained features are not captured or misunderstood. (c) Predictions are not correlated with the given descriptions.

sual dialog [44, 18]. VG aims to locate a target object within an image based on the given language reference expressions by incorporating information from both textual and visual modalities. However, existing VG methods suffer from false-alarm issues, where they assume that the referred object always exists in the image, leading to inaccurate or wrong targets being detected when irrelevant or inaccurate textual expressions are provided, shown in Fig. 1 (a).

Previous works [21, 20, 36] have made significant progress in VG through various techniques. However, the task of cross-modal learning involved in the VG task remains challenging, and current approaches can be broadly divided into two main categories: two-stage methods [15, 30, 4, 46] and one-stage methods [55, 28, 52, 34, 41, 53]. Despite the significant achievements, the VG approaches suffer from some limitations, such as failing to capture the detailed feature representation accurately, resulting in a lack of discrimination between fine-grained objects with reference expressions shown in Fig. 1 (b), and detecting irrelevant or incorrect targets without understanding the whole

context shown in Fig. 1 (c).

To address the above issues, this paper proposes a novel iterative robust visual grounding (IR-VG) approach with masked reference based centerpoint supervision. The approach first constructs two new robust VG datasets and proposes a multi-stage false-alarm sensitive decoder (MFSD) module to handle the case when there is no target object from the textual expression, avoiding generating false alarms. Secondly, a new masked reference based centerpoint supervision (MRCS) module is proposed to capture the fine-grained feature and enhance the localization capacity from the given reference expressions. Finally, an iterative multi-level vision-language fusion (IMVF) module is leveraged to fuse multi-level visual and textual information that are crucial for vision-language understanding.

The contributions of this paper are summarized as follows: firstly, the proposed approach handles the false-alarm issue in VG task for the **first time** by constructing two new robust VG benchmarks and introducing a multi-stage false-alarm sensitive decoder (MFSD) module. Secondly, a new masked reference based centerpoint supervision (MRCS) module is proposed to achieve much more accurate fine-grained feature and better localization capacity from fully visual-textual comprehension. Lastly, the iterative multi-level vision-language fusion (IMVF) module is introduced to comprehensively fuse multi-level visual and textual information for better vision-language understanding and alignment. Extensive experiments on five *regular* VG benchmarks and two newly constructed *robust* VG benchmarks demonstrate the effectiveness of the proposed approach, achieving above **10%** improvement on robust datasets.

## 2. Related Work

**Visual Grounding.** The Visual Grounding task is an important problem in computer vision that aims to localize an object within an image based on a given language reference expression. The existing approaches typically extend the object detection framework, such as YOLOV3 [38], Faster-RCNN [39], RetinaNet [29], CenterNet [10], and DETR [3], by incorporating a visual-linguistic fusion module. These approaches can be categorized into two main categories: *two-stage methods* [15, 30, 4, 46, 57] and *one-stage methods* [55, 28, 52, 34, 41, 53]. Two-stage approaches, including CMN [15], NMTree [30] and RefNMS [4], Two-branch Network [46] and MAttNet [57], utilize an object detector to generate region proposals and then use textual descriptions to select the highest scoring proposal in the second stage. However, this approach can be computationally expensive due to the large number of proposals, and the matching process for each proposal may slow down the inference speed. On the other hand, one-stage approaches [55, 28, 52, 34, 41, 53, 16] directly incorporate the linguistic context into visual features to predict

the object's location, without generating region proposals. Although one-stage approaches are simple and efficient, they typically rely on pointwise feature representations, which may not be flexible enough to achieve a global context understanding from the vision-language information. Recently, *transformer-based* Visual Grounding approaches have gained popularity due to their attention capacity and efficiency. For instance, TransVG [7] captures intra- and inter-modal contexts using transformers in a uniform manner, while VLTVG [51] builds discriminative feature maps and detects the target object through a multi-stage decoder.

**Robustness in Visual Grounding.** Recent studies have explored CNN robustness in various benchmarks [12] [35], and some works have evaluated and improved CNN robustness for practical applications [43] [42] [1] [48]. RefSegformer [48] incorporates negative sentence inputs to handle false-alarm issues in referring segmentation tasks. However, to the best of our knowledge, no existing benchmarks or approaches have explored the robustness of the Visual Grounding task. In practice, existing approaches often fail to generate accurate targets when an irrelevant or inaccurate language expression is given. Therefore, this paper takes a further step by proposing a new iterative *robust* VG framework and building two robust VG datasets to address this research problem. It is important to note that, within the context of this paper, the term "**robust**" refers to the ability of the proposed method to produce accurate results and avoid false-alarm predictions even when provided with irrelevant and incorrect expressions.

**Multi-modal Transformer.** Vision transfomer [3, 9, 24, 61, 60, 26, 62, 25, 27, 6, 50, 32] has a a wide range of application, including detecction, representation learning, and segmentation. Recent works [19, 64, 63, 49, 47] unify different modal inputs and outputs, mainly representation learning, open vocablulary, and large language models. For visual grounding, recent works [48, 56, 23, 7] also adopt multumodal transformer framework, our method belong to this scope. In partilcaur, we pay more attention on the robustness and fine-grained supervision design.

## 3. Method

In this section, we present the architecture of the proposed robust VG pipeline and its components. Fig. 2 illustrates the pipeline. In this section, we present the architecture of the proposed robust VG pipeline and its components. Fig. 2 illustrates the pipeline, where the image and corresponding language description are processed separately to obtain different feature embeddings in two distinct branches.
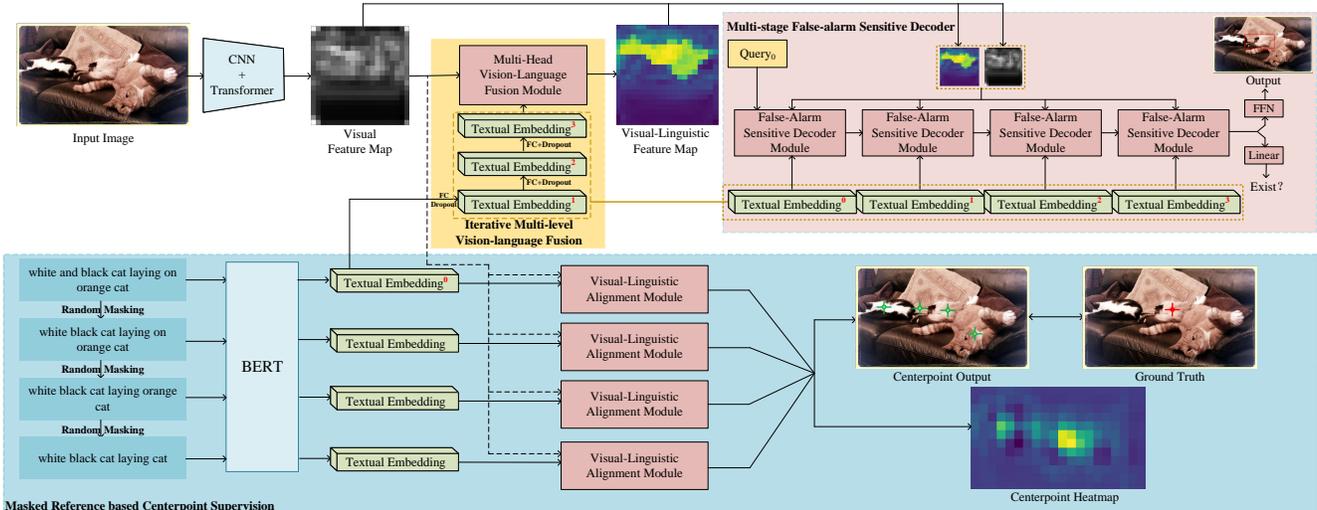
Figure 2. An overview of our proposed IR-VG framework, which comprises Masked Reference based Centerpoint Supervision, Iterative Multi-Level Vision-Language Fusion, and Multi-Stage False-Alarm Sensitive Decoder.

## 3.1. Masked Reference based Centerpoint Supervision

**Motivation.** Existing VG approaches suffer from inadequate visual-linguistic feature representation, insufficient fine-grained feature representation, and poor localization capacity, leading to the detection of irrelevant or inaccurate objects. To address these issues, we propose the masked reference based centerpoint supervision (MRCS) approach, as illustrated in Fig. 2. MRCS comprises three parts: masked reference augmentation, visual-linguistic alignment, and centerpoint supervision. This approach aims to enhance context understanding from the whole image and improve the accuracy of object detection in VG tasks.

**Masked reference augmentation.** As illustrated in the down-left part of Fig. 2, we propose a text augmentation approach to generate diversified textual information given an input language expression. We employ the NLTK [2] tokenization strategy to extract lexical properties for each word, followed by masking one word in the text according to the well-designed rules (shown in the supplementary materials). This masking process is repeated at most three times, achieving one full text and three masked texts in total. BERT [8] is then utilized to generate different textual embeddings for these sentences.

It is important to note that we prioritize the masking of lexical words differently based on their semantic significance. Prepositions, conjunctions, and qualifiers are masked first, as they generally have minimal impact on the sentence's meaning. If these types of words are absent, the module masks auxiliaries, pronouns, and numbers, which can partially affect the sentence's semantics. Finally, the module masks adjectives and verbs, which are critical for the sentence's meaning. If there is only one non-noun word

or only nouns remaining in the sentence, no further masking is performed. More specific rules will be shown in the supplementary materials.

**Visual-linguistic alignment.** The proposed model, illustrated in Fig. 3, incorporates a visual-linguistic alignment module with two consecutive MHA layers. The visual feature map $F_v$ is input as the *Query*, and the textual embeddings are input as *Key* and *Value* to the first MHA. This process produces an enhanced feature map that gathers relevant semantic information from the corresponding linguistic representation. Subsequently, the enhanced feature map undergoes another MHA operation that performs self-attention on the visual features to encode the involved visual contexts. The features from the two MHAs are element-wisely summed in a residual manner for the centerpoint supervision component. The goal of these two MHA operations is to encode the related descriptions into the visual feature and enhance the visual context information from the whole image. The features from the two MHAs are element-wisely summed in a residual manner for the centerpoint supervision component. As shown in Fig. 3, the textual embeddings from the language branch and the visual feature map from the image branch will be input to the visual-linguistic alignment module based on two consecutive multi-head attention (MHA) layers. Specifically, we input the visual feature map $F_v$ as the *Query*, and textual embeddings as *Key* and *Value* into the first MHA layer, where enhanced feature map will be achieved by collecting the relevant semantic information from the corresponding linguistic representation. The enhanced feature will then again be processed through another MHA operator that performs self-attention for the visual features to encode the involved visual contexts. The two consecutive MHA operations try to encode the related
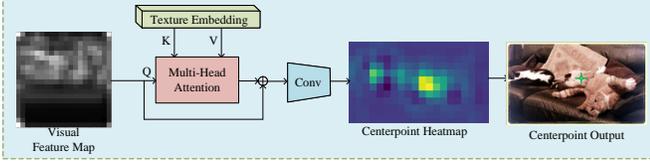
Figure 3. The architecture of visual-linguistic alignment module.
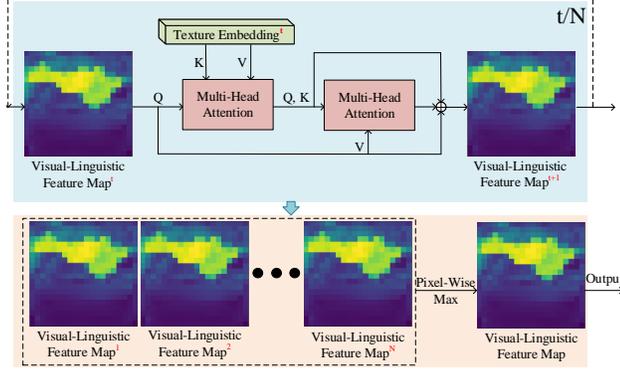

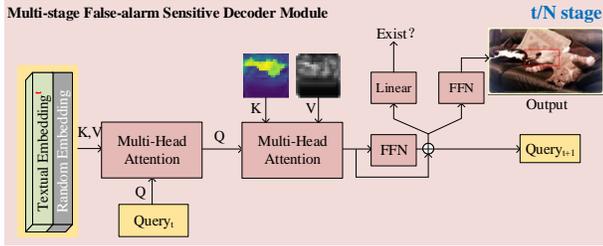
Figure 4. The architecture of IMVF.



Figure 5. The architecture of MFSD.

descriptions into the visual feature and enhance the visual context information from the whole image. The features from the two MHAs will be element-wisely summed together in a residual manner, which will be employed in the keypoint supervision part.

**Centerpoint supervision.** To obtain the final centerpoint heatmaps, the summed feature map obtained from each language expression is processed through two consecutive convolutional layers. Multiple centerpoint heatmaps (one from the full text and three from the masked text) are then fused by performing a *maxpooling* process, with the centerpoint coordinates determined by performing a *argmax* operation on the resulting heatmap. The cross entropy loss is then utilized as the supervision loss between the centerpoint heatmap and the corresponding ground truth, given by $\mathcal{L}_{\text{key}} = \text{CELoss}(y, \hat{y}_i)$, where $\text{CELoss}(\cdot, \cdot)$ is the cross entropy loss, $\hat{y}_i$ is the predicted centerpoints, and $y$ is the centerpoint ground truth that is obtained from the center point of each ground truth box.

## 3.2. Iterative Multi-level Vision-language Fusion

**Motivation.** Through empirical analysis, we have observed that the visual-textual misunderstanding issue arises due to inadequate and poor textual embeddings in the multi-head vision-language fusion module, which incorporates different visual features from various stages. To address this challenge, we propose a multi-level textual feature enhancement (MTFE) module that enhances textual embeddings from low-level to high-level, analogous to the image feature extraction branch. The module extracts multi-level textual information from the entire sentence, resulting in more comprehensive and robust textual embeddings.

**Multi-level textual feature enhancement.** The MTFE module improves textual embedding representation by performing two consecutive fully-connected layers with 768 nodes in each stage. Specifically, as highlighted with yellow color in Fig. 2, the IMVF comprises four stages, and each stage contains an MTFE module. The MTFE module consists of two fully connected layers and a corresponding dropout layer with a 0.1 ratio, aimed at obtaining multi-level textual features that match the multi-level visual features. This enables the model to focus on different key descriptions in the referring expressions and obtain more complete and reliable features for the referred object.

**Iterative multi-level vision-language fusion.** Fig. 4 illustrates the IMVF module, which is based on MHA and consists of four iterative stages. Each stage includes two MHA layers. The first layer uses the visual feature map $F_v \in \mathcal{R}^{C \times H \times W}$ as the *Query* and the textual embeddings $F_l \in \mathcal{R}^{C \times L}$ from the multi-level textual feature enhancement module as the *Key* and *Value*. Multi-head cross-attention enables the comprehensive incorporation of textual information into the visual feature map $F_g \in \mathcal{R}^{C \times H \times W}$. In the second layer, $F_g$ serves as both the *Query* and *Key*, while $F_v$ serves as the *Value*. This self-attention operator allows the model to gather crucial context features for the referred object based on the textual descriptions provided, and the final feature is $F_c \in \mathcal{R}^{C \times H \times W}$. We sum the $F_v$, $F_g$, and $F_c$ element-wisely to obtain the final visual feature map $F_m$. In each iteration, the $i$-th visual feature map $F_m^i$ becomes the initial feature map (i.e., $F_v^{i+1}$). Our experiments include four iterations, and we use element-wise *max* strategy to obtain the final fusion feature $F = \max(F_m^1, F_m^2, F_m^3, F_m^4)$. Actually, other fusion strategies can also be considered. We experimentally find that element-wise summation or product achieves inferior performance than the proposed strategy.

## 3.3. Multi-stage False-alarm Sensitive Decoder

**Motivation.** The current SOTA approaches in VG task assume that the language expressions are precisely matched with the visual image. However, this assumption may not hold in practical applications. Specifically, when an inac-

curate or irrelevant text expression is provided, the existing SOTA VG approaches [7, 51] often generate false-alarm results. To address this issue, we introduce several **robust** VG datasets (described in Sec. 4) and propose a new multi-stage false-alarm sensitive decoder (MFSD) module.

**Multi-stage false-alarm sensitive decoder.** As shown in Fig. 5, the MFSD module consists of several iterative stages, each contains two consecutive multi-head attention (MHA) [45] layers. In the first stage, we randomly initialize a series of learnable queries. To handle the false-alarm case, we introduce a random embedding with the same size as textual embedding from the IMVF module. We concatenate the textual embedding and the random embedding in the batch dimension, termed as *mixture embedding*. For the first MHA layer, the learnable queries serves as *Query*, and the mixture embedding acts as *Key* and *Value*. With this layer, the textual embedding can be more easily attended to the target tokens, thus achieving enhanced textual embedding. For the second MHA layer, the enhanced textual embedding is treated as *Query*, and the visual-linguistic feature map from the IMVF module as well as the visual feature map $F_v$ are employed as *Key* and *Value*. Through the second MHA layer, the textual information can be comprehensively fused with the visual feature map to achieve an enhanced vision-language feature, which is then taken into a feed-forward network (FFN). We fuse the enhanced vision-language feature and the feature from the second MHA in a residual manner, termed as R_feature, which serves as the *Query* in the next iteration. Then, R_feature is taken into two decoupled heads: one for classification to indicate whether there exists false-alarm result, and another for regression to generate the predicted bounding boxes (bbox). Specifically the classification loss $\mathcal{L}_{cls}$ and the regression loss $\mathcal{L}_{reg}$ are defined as,

$$\mathcal{L}_{cls} = \sum_{t=1}^{N} \sum_{i=1}^{K} \text{CELoss}(y^t, \hat{y}_i^t), \tag{1}$$

$$\mathcal{L}_{reg} = \sum_{t=1}^{N} \sum_{i=1}^{K} \lambda_{GIOU} \mathcal{L}_{GIOU}(b^t, \hat{b}_i^t) + \lambda_{L1} \mathcal{L}_{L1}(b^t, \hat{b}_i^t), \tag{2}$$

where $\text{CELoss}(\cdot, \cdot)$, $\mathcal{L}_{GIOU}(\cdot, \cdot)$ and $\mathcal{L}_{L1}(\cdot, \cdot)$ are the cross entropy loss, GIOU loss [40] and L1 loss, respectively. $y^t$ and $\hat{y}_i^t$ denote the ground truth label and predicted result in $t$-th iteration. Similarly, $b^t$ and $\hat{b}_i^t$ denote the ground truth bbox and predicted bbox. $t$ denotes the $t$-th iteration, and $i$ represents the $i$-th bbox. $\lambda_{GIOU}$ and $\lambda_{L1}$ are empirically adjusted, here we set them as 3 and 7 by default for all the following experiments.

### 3.4. Details of Determining False-alarm Detection of the Previous Methods

In the previous methods, we follow the same rules as ours to obtain the false alarm. Firstly, we achieve the top1

scoring box as the final prediction box. Then, we calculate the IOU value with the ground truth box. If the IOU value is greater than 0.5, we consider it a true positive, otherwise, we treat it as a false positive. However, the proposed method differs in that it combines the top1 scoring box and its existing result (exist or non-exist) to achieve the final prediction box. During our experiments, we attempted to add an irrelevant text reference head to some previous networks, such as VLTVG [51] but the results were inferior to their baselines. It may not be fair to compare these results in the paper, thus we do not show these results.

### 3.5. Training Loss

In the training stage, the proposed VG framework is trained end-to-end using the aforementioned losses. The overall loss function for the proposed framework is $\mathcal{L} = \mathcal{L}_{cls} + \lambda_{reg}\mathcal{L}_{reg} + \lambda_{key}\mathcal{L}_{key}$ as follows, where $\mathcal{L}_{cls}$, $\mathcal{L}_{reg}$, and $\mathcal{L}_{key}$ denote the classification loss, regression loss and centerpoint loss, respectively. $\lambda_{reg}$ and $\lambda_{key}$ are introduced to balance the above losses. We empirically set $\lambda_{reg}$ and $\lambda_{key}$ as 2 and 5 by default.

$$\mathcal{L} = \mathcal{L}_{cls} + \lambda_{reg}\mathcal{L}_{reg} + \lambda_{key}\mathcal{L}_{key}, \tag{3}$$

## 4. Experiment

### 4.1. Experimental Settings

**Datasets.** To comprehensively verify the effectiveness of the proposed robust VG approach, we evaluate it on two types of datasets: the regular VG datasets and the robust VG datasets.

**Regular VG datasets.** We evaluate our proposed approach on five regular VG datasets, including the RefCOCO [58], RefCOCO+ [58], RefCOCOg [33], ReferItGame [22], and Flickr30k [37]. The RefCOCO datasets series, including RefCOCO, RefCOCO+, and RefCOCOg, are three commonly used benchmarks for visual grounding, the images used in these datasets are collected from the train2014 set of MSCOCO dataset. Specifically, the RefCOCO dataset contains 19,994 images, 50,000 reference objects, and a total of 142,210 reference expressions. Among them, 120,624 reference expressions are used as the training set, 10,834 as the validation set, 5657 and 5095 expressions for test A and test B, respectively. The RefCOCO+ dataset provides 19,992 images with 49,856 reference objects and 141,564 reference expressions. Similar to RefCOCO, RefCOCO+ is also divided into training, validation, test A, and test B sets, with 120,191, 10,758, 5,726, and 4,889 reference expressions in these datasets. RefCOCOg contains a total of 25,799 images, 49,822 objects, and 95,010 reference expressions. Compared to the first two datasets, most of the expressions in RefCOCOg have longer sentences and more complex statement structures. RefCOCOg contains two sub-datasets, RefCOCOg-google and RefCOCOg-umd.

Since the former dataset does not provide a test set, we mainly use the RefCOCOg-umd dataset. ReferItGame contains 20,000 images, which are collected from the SAIAPR-12 dataset. This dataset has a total of 120,072 reference expressions and is divided into a training set with 54,127 reference expressions, a validation set with 5,842 reference expressions, and a test set with 60,103 reference expressions. Flickr30k contains 31,783 images and 427,000 reference expressions. We divide the training, validation, and test sets using the same ratio as the previous work.

**Robust VG Datasets** We construct two robust VG datasets based on the existing benchmarks RefCOCOg and ReferItGame, termed RefCOCOg_F and ReferItGame_F. The train set of our robust VG datasets contains two parts of data, the first part is the train set of the original dataset, while the second part is a random matching dataset, which destroys the correspondence between the image information and the language descriptions. Specifically, for each target on the image, we select one description that is different from its original one among all the text descriptions in the dataset, thus building a dataset where the image is with irrelevant or inaccurate descriptions. During training, the ratio of these two parts of data is 1:1. The test set of our robust VG datasets also consists of two parts of data, the first part is the test set of the original dataset while the second part is the manually modified robust VG dataset, which requires manual intervention to modify some keywords in the descriptions, thus modifying the semantics of the descriptions and building a more difficult dataset. For instance, we manually modify the expression "The man in white T-shirt is riding a bike" to "The man in blue T-shirt is riding a bike". Specifically, the test set of the RefCOCOg_F dataset contains 2000 pairs of false-alarm data and 9602 pairs of regular data that are from the original RefCOCOg test set. The test set of the ReferItGame_F dataset contains 1000 pairs of false-alarm data and 9000 pairs of regular data that are randomly sampled from the test set of the original ReferItGame dataset.

Specifically, the data combination method of the random matching dataset is to randomly replace the description in each group of data in the training set with a random other description in the dataset to construct false-alarm data. Of course, the description of the same image will not be selected to avoid the existence of the target corresponding to the ran71 dom description on the image. It can be observed that the probability of the existence of the target corresponding to the description on the image is very low for the false alarm data formed by this random selection description method.

We build the manually modified robust VG dataset by manually modifying some keywords in the description. In general, we mainly modify words from the following perspectives. First, modifying key nouns can greatly change



Figure 6. Example of manually modified false-alarm data.

the semantics of words, thus generating false alarm data. For example, modify "Two men on a horse" to "Two men on a car" (as shown in the first row of Fig. 6). Second, modifying key adjectives can also change the description semantics. For example, modify "A man with a bat wearing a red helmet" to "A man with a bat wearing a yellow helmet" (as shown in the second row of Fig. 6). Third, modify words in the text that relate to spatial location can mismatch the original target with the newly generated text. For example, modify "An elephant trainer standing beside an elephant walking down the street" to "An elephant trainer standing far away from an elephant walking down the street" (as shown in the third row of Fig. 6). Fourth, changing the words corresponding to some fine-grained features can generate false-alarm data. For example, modify "A man wearing glasses" to "A man without glasses" (as shown in the fourth row of Fig. 6). Experiments show that our pro95 posed IR-VG is effective for all four types of false alarm data.

**Implementation Details.** Consistent with SOTA approaches such as TransVG [7] and vltvg [51], our proposed method employs ResNet101 [14] as the backbone, augmented with 6 transformer layers in the image feature extraction branch, initialized using weights from DETR [3]. The textual embedding extraction branch is initialized with BERT [8], while the parameters of other components use Xavier scheme [13] initialization. We resize all images to $640 \times 640$ and fill them with black to form a square. We perform experiments using PyTorch, a 3090ti GPU, a batch size of 16, and run training for 90 epochs using the AdamW optimizer with a learning rate of $3 \times 10^{-4}$ and a weight decay of $1 \times 10^{-4}$.

**Evaluation Metrics.** For the *regular* VG datasdet, following previous works [7] [54], we adopt the commonly used top1 accuracy (acc-1) as the evaluation metric. For the *robust* VG dataset, we propose two novel evaluation metrics,

| Method | RefCOCO | | | RefCOCO+ | | | RefCOCOg | | ReferItGame | Flickr30k |
|---|---|---|---|---|---|---|---|---|---|---|
| | val | testA | testB | val | testA | testB | val-u | test-u | test | test |
| CMN [15] | - | 71.03 | 65.77 | - | 54.32 | 47.76 | - | - | 28.33 | - |
| VC [59] | - | 73.33 | 67.44 | - | 58.40 | 53.18 | - | - | 31.13 | - |
| NMTree [30] | 76.41 | 81.21 | 70.09 | 66.46 | 72.02 | 57.52 | 65.87 | 66.44 | - | - |
| Ref-NMS [4] | 80.70 | 84.00 | 76.04 | 68.25 | 73.68 | 59.42 | 70.55 | 70.62 | - | - |
| FAOA [55] | 72.54 | 74.35 | 68.50 | 56.81 | 60.23 | 49.60 | 61.33 | 60.36 | 60.67 | 68.71 |
| LBYLNet [17] | 79.67 | 82.91 | 74.15 | 68.64 | 73.38 | 59.49 | - | - | 67.47 | - |
| TransVG [7] | 81.02 | 82.72 | 78.35 | 64.82 | 70.70 | 56.94 | 68.67 | 67.73 | 70.73 | 79.10 |
| VLTVG [51] | 84.77 | 87.24 | 80.49 | 74.19 | 78.93 | 65.17 | 76.04 | 74.98 | 71.98 | 79.84 |
| IR-VG (Ours) | **86.82** | **88.75** | **82.60** | **76.22** | **80.75** | **67.33** | **77.86** | **76.24** | **74.03** | **81.45** |

Table 1. Comparisons with SOTA visual grounding methods.



Figure 7. Visualization of the MFSD module.

| Methods | RefCOCOg_F | | ReferItGame_F | |
|---|---|---|---|---|
| | $R_{fad}$ | $R_{mix}$ | $R_{fad}$ | $R_{mix}$ |
| CMN [15] | 27.10 | 65.10 | 24.75 | 21.41 |
| VC [59] | 42.45 | 68.85 | 31.03 | 25.69 |
| SSG [5] | 34.15 | 61.25 | 32.44 | 46.43 |
| Ref-NMS [4] | 43.90 | 62.40 | 41.39 | 48.15 |
| ReSC-Large [52] | 37.35 | 60.55 | 32.54 | 59.89 |
| LBYLNet [17] | 45.40 | 63.32 | 45.40 | 60.57 |
| IR-VG (Ours) | **67.32** | **73.61** | **69.44** | **72.03** |

Table 2. Comparisons with SOTA approaches on *robust* VG datasets.

| Methods | | RefCOCO | | | RefCOCO+ | | | RefCOCOg | |
|---|---|---|---|---|---|---|---|---|---|
| I | M | val | testA | testB | val | testA | testB | val-u | test-u |
| - | - | 84.77 | 87.24 | 80.49 | 74.19 | 78.93 | 65.17 | 76.04 | 74.98 |
| ✓ | - | 85.92 | 88.41 | 81.77 | 75.27 | 80.06 | 66.33 | 77.10 | 76.06 |
| - | ✓ | 85.53 | 88.09 | 81.23 | 75.34 | 79.97 | 66.18 | 77.21 | 75.75 |
| ✓ | ✓ | **86.82** | **88.75** | **82.60** | **76.22** | **80.75** | **67.33** | **77.86** | **76.24** |

Table 3. Ablation studies on three benchmarks, "I" and "M" denote IMVF and the MRCS.

i.e., false alarm discovery rate $R_{fad}$ with only false-alarm data, and correct rate among the mixed data $R_{mix}$ with both false-alarm and regular data, which are defined as,

$$R_{fad} = \frac{FA^{acc}}{FA^{all}}, \quad R_{mix} = \frac{FA^{acc} + Regular^{acc}}{FA^{all} + Regular^{all}}, \quad (4)$$

where FA denotes the false-alarm data with irrelevant or inaccurate descriptions, and Regular means the regular data with accurate descriptions. The superscript **acc** and **all** represent the number of accurate predictions and the total number of the data. The detailed dataset descriptions, training loss and other experiment implementation details will be shown in the supplementary materials.

### 4.2. Comparisons with Existing SOTA Methods

As presented in Tab. 1, we evaluate the proposed approach against other SOTA VG methods. Numerically, we improve over the best SOTA approaches by about 2% in all five benchmarks, indicating the effectiveness of our proposed method.

Tab. 2 demonstrates the numerical comparisons on the *robust* VG datasets. Obviously, we improve over the SOTA approaches by a nontrivial margin in competitive benchmarks of RefCOCOg_F and ReferItGame_F. Specifically, on ReferItGame_F dataset, we achieve about 25% and 10% improvement in $R_{fad}$ and $R_{mix}$ metrics, respectively. It is worth noting that TransVG [7] and VLTVG [51] are not included in the comparison because they only provide one predicted bounding box without any extra information to determine whether the target object is a false alarm. As a result, they will definitely generate false-alarm objects when given inaccurate or irrelevant language expressions, which is not a fair comparison.

### 4.3. Ablation Study

**Numerical Component Analysis.** Tab. 3 shows the effectiveness of each component on the *regular* VG datasets. The proposed approach outperforms the baseline by 2.1%
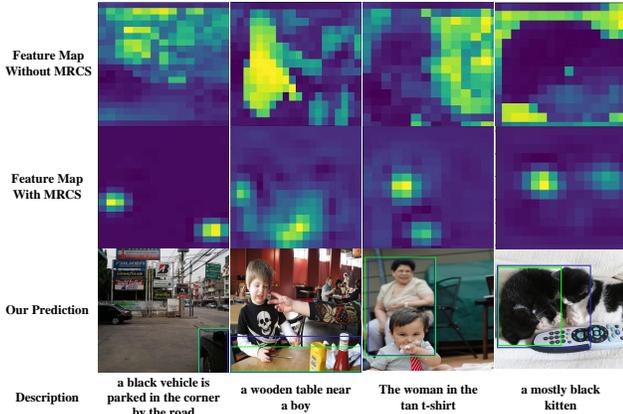
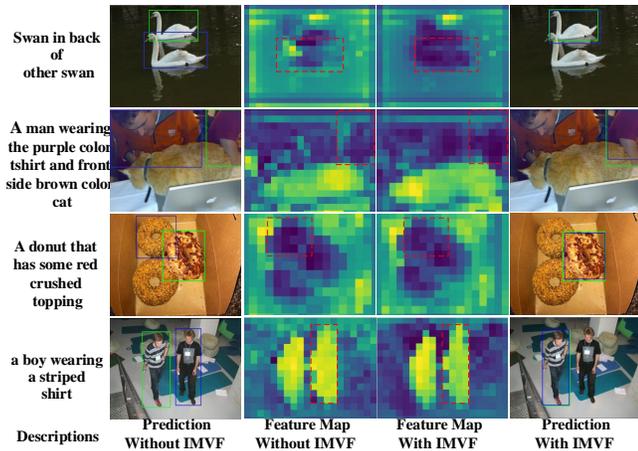Figure 8. Visualization of the visual-linguistic feature map (shown in Fig. 2) with/without MRCS module.



Figure 9. Visualization of the visual-linguistic feature map (shown in Fig. 2) with/without the IMVF module, especially for the red rectangle areas.

top1 accuracy in RefCOCO testB dataset. Specifically, IMVF improves by 1.3% and MRCS improves by 0.7%. Similar conclusions can be drawn from other *regular* VG datasets. Tab. 2 illustrates the effectiveness and robustness of the proposed MFSD module, which achieves a significant improvement on two competitive robust benchmarks. For instance, the MFSD module improves by 25% and 10% compared with existing SOTA approaches in $R_{fad}$ and $R_{mix}$ metrics, respectively.

## 4.4. Rules for masking words in MRCS module.

When we mask lexical words, we prioritize them differently. We first mask prepositions, conjunctions, and qualifiers because they usually do not significantly impact the sentence's meaning. If these types of words are not present, the module then masks auxiliaries, pronouns, and numbers, which can partly affect the sentence's semantics. Finally, the module masks adjectives and verbs, which are critical

| Methods | RefCOCO | | | RefCOCO+ | | | RefCOCOg | |
|---|---|---|---|---|---|---|---|---|
| | val | testA | testB | val | testA | testB | val-u | test-u |
| Baseline | 85.92 | 88.41 | 81.77 | 75.27 | 80.06 | 66.33 | 77.10 | 76.06 |
| Wo masked | 86.57 | 88.52 | 82.26 | 75.84 | 80.41 | 67.02 | 77.57 | 75.93 |
| Ours | **86.82** | **88.75** | **82.60** | **76.22** | **80.75** | **67.33** | **77.86** | **76.24** |

Table 4. Ablation study on multiple masked strategies. "Baseline" denotes the experiment with one full text without centerpoint supervision, "Wo masked" denotes the result with one full text and centerpoint supervision, and "Ours" represents the experiment with MRCS.

for the sentence's meaning. If there is only one non-noun word remaining or only nouns remain in the sentence, no further masking is performed. However, even with this priority order, some important words may still get masked, introducing noise into the training. Nevertheless, we empirically demonstrate that the language comprehension improvement from masking operations outweighs the negative effects of introducing noise (shown in Tab. 4). In all datasets, the number of words exceeds 3, and through three masking operations, we find that the majority of the masked words are prepositions, conjunctions, and qualifiers. Therefore, in most cases, this operation will not affect the meaning of the sentence.

**Qualitative Component Analysis.** *Qualitative analysis of MRCS.* Fig. 7 illustrates the visualization of prediction results with or without the MFSD module on the robust VG datasets. It shows that the MFSD module enables the model to efficiently identify the presence or absence of targets described in the text on the image. The first row of the figure shows the false alarm data generated by the key nouns in the description being changed, the second row shows the false alarm data generated by the modification of key adjectives (e.g., color). The third line of the figure shows the spatial location relations in the description being modified and the fourth row of the figure shows the fine-grained features in the description being modified. Our MFSD module can effectively identify the false alarm data generated by all the above modification methods. *Qualitative analysis of MRCS.* Fig. 8 presents the visual-linguistic feature map with or without MRCS module. We intuitively observe that the MRCS enables the feature map to attend more accurately to the target object's location, and generates a more precise foreground map. To avoid interactions from IMVF module, we conduct this experiment only with MRCS module and MFSD module. *Qualitative analysis of IMVF.* Fig. 9 illustrates the visual-linguistic feature map with or without IMVF module. The figure indicates that the IMVF module reduces interference and allows the model to concentrate more on target by better understanding visual and textual information. To ensure fairness, we performed this experiment only with IMVF module and MFSD module.

# 5. Conclusions

Our work introduces the IR-VG framework, which comprises IMVF, MRCS, and MFSD. It outperforms existing approaches in terms of context features, fine-grained features, and localization accuracy while addressing robustness issues when faced with irrelevant or inaccurate reference expressions. Our experiments demonstrate the effectiveness of each module, achieving new SOTA performance.

**Limitation and future work.** Notably, IR-VG builds a new research direction for robust VG. Future work includes developing a more elegant framework to handle false alarms. In addition, we will explore the false-alarm problems with irrelevant expression for some foundation models (e.g. Grounding DINO [31]).

# References

[1] Said Fahri Altindis, Yusuf Dalva, and Aysegul Dundar. Benchmarking the robustness of instance segmentation models. *CoRR*, 2021. 2

[2] Steven Bird, Ewan Klein, and Edward Loper. *Natural Language Processing with Python*. O'Reilly, 2009. 3

[3] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020. 2, 6

[4] Long Chen, Wenbo Ma, Jun Xiao, and et al. Ref-nms: Breaking proposal bottlenecks in two-stage referring expression grounding. In *AAAI*, 2021. 1, 2, 7

[5] Xinpeng Chen, Lin Ma, Jingyuan Chen, and et al. Real-time referring expression comprehension by single-stage grounding network. *CoRR*, 2018. 7

[6] Bowen Cheng, Alexander G. Schwing, and Alexander Kirillov. Per-pixel classification is not all you need for semantic segmentation. 2021. 2

[7] Jiajun Deng, Zhengyuan Yang, Tianlang Chen, and et al. Transvg: End-to-end visual grounding with transformers. In *CVPR*, 2021. 2, 5, 6, 7

[8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*, 2019. 3, 6

[9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. 2021. 2

[10] Kaiwen Duan, Song Bai, Lingxi Xie, Honggang Qi, Qingming Huang, and Qi Tian. Centernet: Keypoint triplets for object detection. In *ICCV*, 2019. 2

[11] Zipeng Fu, Ashish Kumar, Ananye Agarwal, and et al. Coupling vision and proprioception for navigation of legged robots. In *CVPR*, 2022. 1

[12] Robert Geirhos, Carlos R. Medina Temme, Jonas Rauber, and et al. Generalisation in humans and deep neural networks. In *NeurIPS*, 2018. 2

[13] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In Yee Whye Teh and D. Mike Titterington, editors, *AISTATS*, 2010. 6

[14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 6

[15] Ronghang Hu, Marcus Rohrbach, Jacob Andreas, and et al. Modeling relationships in referential expressions with compositional modular networks. In *CVPR*, 2017. 1, 2, 7

[16] Ronghang Hu, Marcus Rohrbach, and Trevor Darrell. Segmentation from natural language expressions. In *ECCV*, 2016. 2

[17] Binbin Huang, Dongze Lian, Weixin Luo, and et al. Look before you leap: Learning landmark features for one-stage visual grounding. In *CVPR*, 2021. 7

[18] Xiaoze Jiang, Jing Yu, Yajing Sun, and et al. DAM: deliberation, abandon and memory networks for generating detailed and non-repetitive responses in visual dialogue. 2020. 1

[19] Aishwarya Kamath, Mannat Singh, Yann LeCun, and et al. MDETR - modulated detection for end-to-end multi-modal understanding. In *CVPR*, 2021. 2

[20] Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *CVPR*, 2015. 1

[21] Andrej Karpathy, Armand Joulin, and Li Fei-Fei. Deep fragment embeddings for bidirectional image sentence mapping. 2014. 1

[22] Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara L. Berg. Referitgame: Referring to objects in photographs of natural scenes. In *EMNLP*, 2014. 5

[23] Namyup Kim, Dongwon Kim, Cuiling Lan, Wenjun Zeng, and Suha Kwak. Restr: Convolution-free referring image segmentation using transformers. 2022. 2

[24] Xiangtai Li, Henghui Ding, Wenwei Zhang, Haobo Yuan, Jiangmiao Pang, Guangliang Cheng, Kai Chen, Ziwei Liu, and Chen Change Loy. Transformer-based visual segmentation: A survey. *arXiv preprint arXiv:2304.09854*, 2023. 2

[25] Xiangtai Li, Shilin Xu, Yibo Yang, Haobo Yuan, Guangliang Cheng, Yunhai Tong, Zhouchen Lin, Ming-Hsuan Yang, and Dacheng Tao. Panopticpartformer++: A unified and decoupled view for panoptic part segmentation. *arXiv preprint arXiv:2301.00954*, 2023. 2

[26] Xiangtai Li, Haobo Yuan, Wenwei Zhang, Guangliang Cheng, Jiangmiao Pang, and Chen Change Loy. Tube-link: A flexible cross tube baseline for universal video segmentation. *ICCV*, 2023. 2

[27] Xiangtai Li, Wenwei Zhang, Jiangmiao Pang, Kai Chen, Guangliang Cheng, Yunhai Tong, and Chen Change Loy. Video k-net: A simple, strong, and unified baseline for video segmentation. In *CVPR*, 2022. 2

[28] 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

[29] Tsung-Yi Lin, Priya Goyal, Ross B. Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *ICCV*, 2017. 2

[30] Daqing Liu, Hanwang Zhang, Zheng-Jun Zha, and et al. Learning to assemble neural module tree networks for visual grounding. In *CVPR*, 2019. 1, 2, 7

[31] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023. 9

[32] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, 2021. 2

[33] Junhua Mao, Jonathan Huang, Alexander Toshev, Oana Camburu, Alan L. Yuille, and Kevin Murphy. Generation and comprehension of unambiguous object descriptions. In *CVPR*, 2016. 5

[34] Yue Ming, Nannan Hu, Chunxiao Fan, and et al. Visuals to text: A comprehensive review on automatic image captioning. *IEEE/CAA Journal of Automatica Sinica*, 2022. 1, 2

[35] Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and et al. Exploring generalization in deep learning. In *NeurIPS*, 2017. 2

[36] Ahmad Ostovar, Suna Bensch, and Thomas Hellström. Natural language guided object retrieval in images. *Acta Informatica*, 2021. 1

[37] 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. *CVPR*, 2017. 5

[38] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *CoRR*, 2018. 2

[39] Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. Faster R-CNN: towards real-time object detection with region proposal networks. *PAMI*, 2017. 2

[40] Hamid Rezatofighi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian D. Reid, and Silvio Savarese. Generalized intersection over union: A metric and a loss for bounding box regression. In *CVPR*, 2019. 5

[41] Arka Sadhu, Kan Chen, and Ram Nevatia. Zero-shot grounding of objects from natural language queries. In *ICCV*, 2019. 1, 2

[42] Christos Sakaridis, Dengxin Dai, and Luc Van Gool. Semantic foggy scene understanding with synthetic data. *IJCV*, 2018. 2

[43] Christos Sakaridis, Dengxin Dai, and Luc Van Gool. Guided curriculum model adaptation and uncertainty-aware evaluation for semantic nighttime image segmentation. In *CVPR*, 2019. 2

[44] Kaili Sun, Chi Guo, Huyin Zhang, and et al. HVLM: exploring human-like visual cognition and language-memory network for visual dialog. *Information Processing & Management*, 2022. 1

[45] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 5

[46] Liwei Wang, Yin Li, Jing Huang, and Svetlana Lazebnik. Learning two-branch neural networks for image-text matching tasks. *PAMI*, 2019. 1, 2

[47] Jianzong Wu, Xiangtai Li, Henghui Ding, Xia Li, Guangliang Cheng, Yunhai Tong, and Chen Change Loy. Betrayed by captions: Joint caption grounding and generation for open vocabulary instance segmentation. *ICCV*, 2023. 2

[48] Jianzong Wu, Xiangtai Li, Xia Li, Henghui Ding, Yunhai Tong, and Dacheng Tao. Towards robust referring image segmentation. *CoRR*, 2022. 2

[49] Jianzong Wu, Xiangtai Li, Shilin Xu, Haobo Yuan, Henghui Ding, Yibo Yang, Xia Li, Jiangning Zhang, Yunhai Tong, Xudong Jiang, Bernard Ghanem, and Dacheng Tao. Towards open vocabulary learning: A survey. *arXiv pre-print*, 2023. 2

[50] Shilin Xu, Xiangtai Li, Jingbo Wang, Guangliang Cheng, Yunhai Tong, and Dacheng Tao. Fashionformer: A simple, effective and unified baseline for human fashion segmentation and recognition. *ECCV*, 2022. 2

[51] Li Yang, Yan Xu, Chunfeng Yuan, and et al. Improving visual grounding with visual-linguistic verification and iterative reasoning. In *CVPR*, 2022. 2, 5, 6, 7

[52] Zhengyuan Yang, Tianlang Chen, Liwei Wang, and et al. Improving one-stage visual grounding by recursive sub-query construction. In *ECCV*, 2020. 1, 2, 7

[53] Zhengyuan Yang, Tianlang Chen, Liwei Wang, and et al. Improving one-stage visual grounding by recursive sub-query construction. In *ECCV*, 2020. 1, 2

[54] Zhengyuan Yang, Tianlang Chen, Liwei Wang, and Jiebo Luo. Improving one-stage visual grounding by recursive sub-query construction. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *ECCV*, 2020. 6

[55] Zhengyuan Yang, Boqing Gong, Liwei Wang, and et al. A fast and accurate one-stage approach to visual grounding. In *CVPR*, 2019. 1, 2, 7

[56] Zhao Yang, Jiaqi Wang, Yansong Tang, Kai Chen, Hengshuang Zhao, and Philip HS Torr. Lavt: Language-aware vision transformer for referring image segmentation. 2022. 2

[57] Licheng Yu, Zhe Lin, Xiaohui Shen, Jimei Yang, Xin Lu, Mohit Bansal, and Tamara L. Berg. Mattnet: Modular attention network for referring expression comprehension. In *CVPR*, 2018. 2

[58] Licheng Yu, Patrick Poirson, Shan Yang, Alexander C. Berg, and Tamara L. Berg. Modeling context in referring expressions. In *ECCV*, 2016. 5

[59] Hanwang Zhang, Yulei Niu, and Shih-Fu Chang. Grounding referring expressions in images by variational context. In *CVPR*, 2018. 7

[60] Jiangning Zhang, Xiangtai Li, Jian Li, Liang Liu, Zhucun Xue, Boshen Zhang, Zhengkai Jiang, Tianxin Huang, Yabiao Wang, and Chengjie Wang. Rethinking mobile block for efficient neural models. *ICCV*, 2023. 2

[61] Jiangning Zhang, Xiangtai Li, Yabiao Wang, Chengjie Wang, Yibo Yang, Yong Liu, and Dacheng Tao. Eatformer: Improving vision transformer inspired by evolutionary algorithm. *arXiv preprint arXiv:2206.09325*, 2022. 2

[62] Qianyu Zhou, Xiangtai Li, Lu He, Yibo Yang, Guangliang Cheng, Yunhai Tong, Lizhuang Ma, and Dacheng Tao. Transvod: end-to-end video object detection with spatial-temporal transformers. *PAMI*, 2023. 2

[63] Jinguo Zhu, Xizhou Zhu, Wenhai Wang, Xiaohua Wang, Hongsheng Li, Xiaogang Wang, and Jifeng Dai. Uni-perceiver-moe: Learning sparse generalist models with conditional moes. *arXiv preprint arXiv:2206.04674*, 2022. 2

[64] Xizhou Zhu, Jinguo Zhu, Hao Li, Xiaoshi Wu, Xiaogang Wang, Hongsheng Li, Xiaohua Wang, and Jifeng Dai. Uni-perceiver: Pre-training unified architecture for generic perception for zero-shot and few-shot tasks. *arXiv preprint arXiv:2112.01522*, 2021. 2