

# YORO - Lightweight End to End Visual Grounding

Chih-Hui Ho<sup>2\*</sup>, Srikar Appalaraju<sup>1\*\*</sup>, Bhavan Jasani<sup>1</sup>, R. Manmatha<sup>1</sup>, and  
Nuno Vasconcelos<sup>2</sup>

<sup>1</sup> AWS AI Labs {srikara,bjasani,manmatha}@amazon.com

<sup>2</sup> UC San Diego {chh279,nuno}@ucsd.edu

**Abstract.** We present YORO - a multi-modal transformer encoder-only architecture for the Visual Grounding (VG) task. This task involves localizing, in an image, an object referred via natural language. Unlike the recent trend in the literature of using multi-stage approaches that sacrifice speed for accuracy, YORO seeks a better trade-off between speed and accuracy by embracing a single-stage design, without CNN backbone. YORO consumes natural language queries, image patches, and learnable detection tokens and predicts coordinates of the referred object, using a single transformer encoder. To assist the alignment between text and visual objects, a novel patch-text alignment loss is proposed. Extensive experiments are conducted on 5 different datasets with ablations on architecture design choices. YORO is shown to support real-time inference and outperform all approaches in this class (single-stage methods) by large margins. It is also the fastest VG model and achieves the best speed/accuracy trade-off in the literature.

## 1 Introduction

There has been significant recent interest in Visual Grounding (VG) [79, 81, 49, 8, 85, 72, 77, 43, 88, 40, 83, 75, 74]. This aims to localize, in an image, an object referred to by natural language, using a text query (see Fig. 1(a)). VG is potentially useful for many applications, ranging from cloud-based image retrieval from large corpora to resource constrained problems on devices of limited computational capacity. For example, in Fig. 1 (a), a robot must use its limited computing resources to resolve a natural language query and locate the desired object. These applications currently pose an extreme challenge to VG, by compounding the difficulty of the task itself with the need to solve it efficiently. When this is the case, computational efficiency becomes an unavoidable requirement of architecture design. Several examples exist in the vision literature, where branches of computationally efficient architectures have evolved for problems like object recognition [23] or detection [58], among others. The goal is not necessarily to develop the most powerful method in the literature, but to find

---

\* Work done during internship at Amazon.

\*\* Corresponding author.

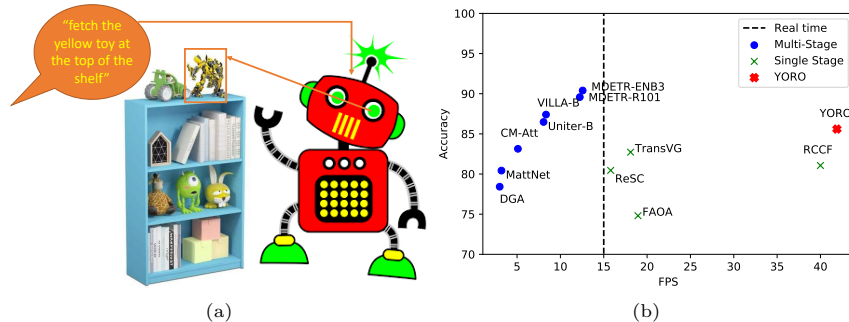


Fig. 1: (a) Example of VG with limited computational resources. (b) **Accuracy-Speed plot**: YORO surpasses real time FPS 15 and achieves better trade-off in the literature. the best trade-off [26] between task performance and some variable that reflects computation, e.g. speed or latency. This has led to the introduction of architectures, such as the MobileNet family [23, 65, 22] for recognition or the YOLO family [58, 59, 60, 1] for object detection, that have become popular despite their weaker than state of the art accuracy. Given the obvious relationship between object detection and VG, it is not surprising that similar trends are starting to emerge for the latter. Like object detectors, VG methods can be broadly grouped in two categories: single [77, 38, 4, 48, 76, 9] and multi-stage [5, 79, 43, 67, 18] methods. Multi-stage methods typically rely on a pre-trained visual backbone (such as the FasterRCNN [61]) as a visual encoder module. This simplifies the VG problem, reducing it to the selection of the object proposals that best match the query text. Like multi-stage object detectors, they are more accurate than single stage methods, but usually too complex for the applications of Fig. 1 (a). In fact, the run-time complexity of many of these VG systems is lower-bounded by the run-time computation of the FasterRCNN [61] proposal stage, which is usually considered too expensive even for object detection, in domains like robotics. In these domains, single-stage detectors such as YOLO [58, 59, 60, 1] are the architecture of choice. Unsurprisingly, the YOLO family has been the backbone of choice for various single-stage VG systems [4, 77]. However, these and the other CNN-based single stage VG systems [77, 38, 4, 48, 76] previously proposed in the literature lack the powerful attention mechanisms that enable recent transformer architectures to excel at multi-modal data fusion [9]. This usually results in weak accuracy.

In this work, we seek to endow single stage VG methods with transformer-style attention, so as to enable a better trade-off between accuracy and speed than currently available methods. For this, we propose an efficient and end-to-end differentiable VG architecture, denoted as *You Only Refer Once* (YORO). YORO is inspired by a transformer architecture, ViT [13], that has achieved success for image classification. ViT has also recently been extended to a transformer architecture, ViLT [32], for vision-language tasks such as visual reasoning [68] or visual question answering [19]. We investigate the more challenging extension of ViLT for tasks that require object localization, such as VG. The main challenge is how to establish explicit connections between words and bounding boxes without

the addition of an object detector. Inspired by recent approaches to the design of transformers for object detection [29, 2, 16], we augment the ViLT architecture with a set of detection tokens that enable this localization. This makes YORO an encoder-only multi-modal transformer architecture that consumes three data streams: referred text tokens, an image in the form of flattened patches, and learnable detection tokens. YORO then uses self-attention to fuse information from all these streams in order to localize the referred object.

Similarly to transformer-based detectors [2, 16], YORO is trained by Hungarian matching between the localization hypotheses derived from its tokens and the ground-truth bounding boxes. This combines a bounding box regression loss that minimizes the difference between predicted and ground truth boxes and a classification loss that associates the referred text to the bounding box. The alignment between language and vision representations is further strengthened by an object-text alignment loss that operates in feature space, encouraging the visual features of the predicted object to match the text features of the corresponding text. However, this must be done carefully. In the early stages of training, when object predictions tend to be incorrect, this object-text alignment loss can hamper learning. To both address this issue and encourage more fine-grain alignment, we propose a novel patch-text alignment loss that aligns features of input patches to their corresponding text. Extensive experiments on five datasets show that YORO establishes the new SOTA for single stage VG models. Overall, as illustrated in Fig. 1(b), YORO (red) achieves the best accuracy/speed trade-off among the methods in the literature. It has substantially better accuracy than previous single stage VG methods (green) and is substantially faster than multi-stage VG models (blue). Ablation studies compare different variants of the proposed model and demonstrate the effectiveness of the proposed patch-text alignment loss. With this we summarize our contributions:

- A novel multi-modal visual grounding architecture YORO is proposed using an encoder-only transformer (Sec. 3). YORO is a simple single-stage design and is end-to-end trainable.
- A novel patch-text alignment loss is introduced for enabling fine-grained visual language alignment (Sec. 4).
- YORO achieves the best accuracy/speed trade-off in the literature, outperforming existing real-time VG methods on five datasets (Sec. 5).

## 2 Related work

Prior work in visual grounding (VG) may be mainly classified into multi-stage and single-stage approaches based on the design of the visual branch.

**Multi-stage VG approaches** [5, 18, 46, 47, 14, 7, 82, 79] usually leverage a region proposal network (RPN) [61] to handle the selection of candidate object locations. In a first-stage, the RPN proposes a set of object candidates. In a second-stage, the VG model selects the region proposal that best matches the query text. This reduces VG to a retrieval problem [24, 25, 63, 17], where the

model only needs to search within the proposed object candidates. Similarly to two-stage object detectors [42, 28], these methods are powerful for VG. However, like two-stage detectors, they tend to be computationally heavy, making them impractical for real-time operation on complexity constrained devices, as illustrated in Fig. 1. When compared to the object detection literature [42, 28], the optimization of the trade-off between speed and performance has received much less attention in VG. On the contrary, recent VG methods [7, 29] pursue higher accuracy by adopting more powerful language models [55, 44] and encoder-decoders transformer architectures [55, 70, 53, 84, 86, 56, 36] with more parameters and slower speeds. Instead, YORO pursues a better trade-off between accuracy and speed. In this context, it is not clear that a multi-stage architecture is the best solution. Since there is no way to recover objects missed by the RPN, strong performance can require a large number of proposals [83, 50], which is inherently expensive. In fact, the use of a preliminary object selection prevents the exploitation of contextual clues that may be available in the query text (e.g. the component “on the top of the shelf” of the query of Fig. 1) and could be critical for solving the task with less proposals. The integration of language and vision is, after all, the difference between object detection and VG. This suggests that single stage architectures, like YORO, could be more effective when computation has a cost. Our experimental results support this hypothesis.

**Single-stage VG approaches** [77, 38, 4, 48, 76, 9] achieve greater inference speeds by discarding the proposal generation stage. Most of these methods are based on CNN object detectors, namely the single-stage detector YOLOV3 [60]. This is, for example, adopted to extract visual features by SSG [4] and FAOA [77]. These features are then fused with text features to regress the bounding boxes. ReSC [76] improves FAOA [77] with a recursive sub-query construction module. MCN [48] further extends YOLOV3 [60] for the joint solution of VG and segmentation, using multitask losses. These models lack the powerful attention mechanisms of the transformer architecture, which are critical for the integration of information both spatially, across image regions, and modally, across the vision and language information streams. When compared to them, YORO has the benefit of leveraging the transformer to implement this powerful integration of information. This enables substantially higher accuracy than previous single-stage approaches. Recently, TransVG [9] has also adopted a transformer model [70] to fuse language and vision features. While TransVG does not use a two-stage object detector, it uses two stages of transformers. The vision and language streams are first processed by independent transformers, whose outputs are then fed to a multi-modal transformer that integrates the two modalities. This makes it substantially slower than YORO, which performs all operations with a single transformer, and could give rise to loss of information needed to reason jointly about vision and language, as discussed above for the RPN. The fact that TransVG is both slower and less accurate than YOLO suggests that there is no benefit to the two transformer stage architecture.

**Vision Transformer.** YORO is based on the Vision transformer (ViT), introduced in [13] for image classification by direct consumption of image patches,



i.e. without a preliminary CNN backbone. Beyond classification [13, 54, 69], the vision transformer has been applied to detection [2, 87, 16] and segmentation [57, 66], as discussed in the recent surveys of [31, 20]. Recently, ViLT [32] has extended the ViT for various visual-language tasks that do not require object localization, such as visual question answering. The extension to tasks, such as VG, that require object localization is non-trivial. This is due to the need to establish explicit connections between words and bounding boxes, which requires much more fine-grained understanding of both language and images. YORO addresses this limitation by the addition of learnable detection tokens and Hungarian matching between token-based object predictions and ground truth bounding boxes, a component of modern transformer-based object detectors [2, 87, 16]. To the best of our knowledge, it is the first architecture to perform VG tasks with an encoder-only transformer without a visual backbone. This is critical for speed. YORO is shown to achieve the best trade-off between speed and accuracy in the VG literature, as illustrated in Fig. 1 (b).

### 3 YORO Architecture

The YORO architecture is summarized in Fig. 2. The input query sentence and image are first converted into text and visual embeddings respectively. These embeddings are augmented by a set of learnable detection tokens and fed into a transformer encoder. Within the transformer, learning happens via self-attention [70]. The transformer predicted features corresponding to the detection tokens are finally fed into several heads, implemented with a multi-layer perceptron, which produce class and bounding box predictions. The model is trained end-to-end, using a combination of losses. Each model component is discussed next.

#### 3.1 Multi-Modal Inputs

YORO consumes multi-modal inputs: language, vision and detection tokens.

**The Language Modality** is depicted in the blue block of Fig. 2. The referred text is converted into sequences of tokens with a BERT tokenizer [10]. Assume the sequence has length  $m$  and a dictionary size of  $v$  for the pre-trained tokenizer. Each text token  $t_i \in \mathbb{R}^v$  in the sequence  $\{t_1; \dots; t_m\}$  is then projected into a  $d$  dimensional feature vector  $f_i^l = \mathcal{W}^l t_i$ , where  $\mathcal{W}^l \in \mathbb{R}^{d \times v}$  is a projection layer. The  $m$  projected token vectors are then augmented with a text classification token  $f_{cls}^l \in \mathbb{R}^d$ . Similar to [32], each of these tokens is summed with a text positional embedding  $f_{pos}^l \in \mathbb{R}^{(m+1) \times d}$  and a modality type embedding  $f_{type}^l \in \mathbb{R}^{(m+1) \times d}$ . The positional embedding specifies the location of each token, while the modality type embedding tells apart the text input from visual input. The overall text input embedding may be written as

$$f^l = [f_{cls}^l; f_1^l; \dots; f_m^l] + f_{pos}^l + f_{type}^l. \quad (1)$$

**The Vision Modality** is shown in the green block of Fig. 2. The input image  $I \in \mathbb{R}^{H \times W \times 3}$  is partitioned into  $n$  visual patches  $\{p_1, \dots, p_n\}$  of size  $s$ , where  $p_i \in \mathbb{R}^{s^2 \times 3}$  and  $n = \frac{HW}{s^2}$ . Each visual patch is then projected, with a linear layer

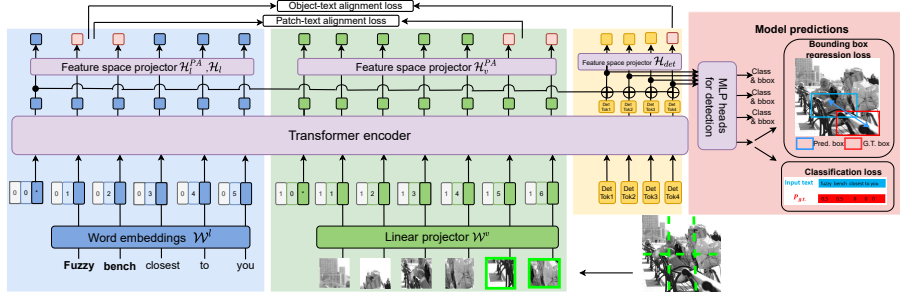


Fig. 2: **YORO architecture**. Blue, green, yellow and pink blocks represent language, vision, detection and prediction branches respectively. YORO does **not** use a large pre-trained visual backbone. Input image is divided into patches. Those of  $\text{IOU} \geq 0.5$  with ground truth box (red) are marked in light green.

of weight matrix  $\mathcal{W}^v \in \mathbb{R}^{d \times (s^2 \times 3)}$ , into a  $d$  dimensional visual patch feature  $f_i^v = \mathcal{W}^v p_i$ . Similar to the language input, the  $n$  visual patch features are augmented with a visual classification token  $f_{cls}^v \in \mathbb{R}^d$  and summed with a visual positional embedding  $f_{pos}^v \in \mathbb{R}^{(n+1) \times d}$  and a modality type embedding  $f_{type}^v \in \mathbb{R}^{(n+1) \times d}$ , to produce an overall visual input embedding

$$f^v = [f_{cls}^v; f_1^v; \dots; f_n^v] + f_{pos}^v + f_{type}^v. \quad (2)$$

**Learnable Detection Tokens.** The vision and language models above are similar to those implemented by ViT [13] and ViLT [32]. These models were proposed for tasks, such as image classification or visual question-answering that do not require object localization. In this case, the use of holistic classifications tokens  $f_{cls}^l, f_{cls}^v$  is sufficient to encode the image identity. However, our experiments (see ablations in Sec. 5.3) show that this is not the case for VG, which depends critically on object localization. In the object detection literature, this problem is addressed by introducing learnable detection tokens, which are individually associated with bounding boxes [2]. YORO leverages such tokens, as shown in the yellow block of Fig. 2. The text  $f^l$  and visual  $f^v$  embeddings are complemented by  $q$  learnable detection tokens  $f^{det} = [f_1^{det}; \dots; f_q^{det}]$ . Each token represents an object and different tokens may represent objects of different sizes.

### 3.2 Transformer Encoder

YORO transformer encoder consists of  $D$  stacked layers, each including a multi-headed self-attention (MSA) network and a feed forward network (FFN). Denoting  $x_d$  as the input of layer  $d$ , the transformer output is computed as follows:

$$x_0 = [f^l; f^v; f^{det}] \quad (3)$$

$$x'_d = x_{d-1} + \text{MSA}(\text{LN}(x_{d-1})), \quad d = 1 \dots D \quad (4)$$

$$x_d = x'_d + \text{FFN}(\text{LN}(x'_d)), \quad d = 1 \dots D \quad (5)$$

where LN denotes layer normalization [70].  $D$  is set to 12 in our experiments.

### 3.3 Multi-Modal Transformer Outputs

The transformer input  $[f^l; f^v; f^{det}]$  is a sequence of length  $n + m + q + 2$ . The output is a sequence  $[o^l; o^v; o^{det}]$  of the same length, where  $o^l = [o_{cls}^l; o_1^l; \dots; o_m^l]$ ,  $o^v = [o_{cls}^v; o_1^v; \dots; o_n^v]$  and  $o^{det} = [o_1^{det}; \dots; o_q^{det}]$ , are the outputs of the language, vision and detection branches respectively. The transformed text class token  $o_{cls}^l$  is added to the transformed detection tokens  $o_i^{det}$  before being passed to the detection heads. This strengthens the dependence of the bounding box predictions on the input text, beyond what would be possible by the use of attention alone.

### 3.4 Feature Projector and Detection Heads

To optimize YORO with the proposed losses, several modules are added to the transformer output, including a detection head for the bounding box and classification loss (red block in Fig. 2), and a projection head to enable the object-text alignment and the patch-text alignment loss discussed in Sec. 4. All these heads are stacked fully connected layers. We next discuss the details of all losses.

## 4 Training Losses

YORO is trained with 4 different losses: (1) a bounding box regression loss to regress bounding box coordinates, (2) a classification loss to classify which text tokens the bounding box is associated to, (3) an object-text alignment loss that aligns detection token features with corresponding text features and (4) a novel patch-text alignment loss that aligns image patches and text tokens. During training, Hungarian matching [34] computes the optimal bipartite matching between predicted and ground truth boxes. During inference, the bounding box of highest classification score is selected.

**Bounding Box Regression Loss.** Two losses are used for bounding box regression: the  $L_1$  and the generalized IoU loss [62]  $L_{giou}$ . These losses are applied to ground truth bounding box  $B_{g.t.} \in \mathbb{R}^4$  and predicted box  $B_{pred} \in \mathbb{R}^4$  pairs. This results in bounding box regression loss

$$L^{bbox} = \lambda_1 L_1(B_{pred}, B_{g.t.}) + \lambda_2 L_{giou}(B_{pred}, B_{g.t.}), \quad (6)$$

where  $\lambda_1$  and  $\lambda_2$  are experimentally set to 2 and 5 respectively.

**Classification Loss.** Unlike the classification losses of object detection, where a single object class is predicted, YORO predicts the set of language tokens to be associated with a bounding box. This is needed to encourage the association of the text phrase with the corresponding objects. For example, in Fig. 2 the red box in the image is associated with the words “fuzzy” and “bench” with probabilities 0.5 and 0.5, respectively. This is implemented by generating predictions with a softmax output of dimension equal to the maximum token length and using a soft cross entropy loss  $L_{sce}$  based on the ground-truth text/object assignments. This is the classification loss  $L^{cls} = L_{sce}(P_{pred}, P_{g.t.})$ , where  $P_{pred}$  is the predicted probability distribution and  $P_{g.t.}$  the ground truth distribution.

**Object-Text Alignment Loss.** Inspired by recent uses of region-text alignment for fusing words and image regions [29, 5, 37, 11, 35, 3], YORO further aligns the language transformer output  $o^l$  and the detection transformer output  $o^{det}$  in a suitable feature space.  $o^l$  and  $o^{det}$  are first mapped to a common feature space  $\mathcal{F}$ , with projection heads  $\mathcal{H}_l$  and  $\mathcal{H}_{det}$ , which implement the mappings  $\hat{o}^l = \mathcal{H}_l(o^l)$  and  $\hat{o}^{det} = \mathcal{H}_{det}(o^{det})$ , where  $\mathcal{H}_{det}$ ,  $\mathcal{H}_{det}$  are linear projections. Consider an image with  $g$  ground truth bounding boxes and let  $\mathcal{T}_i \subseteq \{j : 1 \leq j \leq m\}$ , where  $m$  is the length of the language token sequence, be the subset of language token indices that a projected detection feature  $\hat{o}_i^{det}$  should be aligned to. This alignment is encouraged with the object-to-text alignment (OTA) loss

$$L^{OTA} = \sum_{i=1}^g \frac{1}{|\mathcal{T}_i|} \sum_{j \in \mathcal{T}_i} -\log \left( \frac{\exp(\hat{o}_i^{det\top} \hat{o}_j^l / \tau)}{\sum_{k=1}^m \exp(\hat{o}_i^{det\top} \hat{o}_k^l / \tau)} \right), \quad (7)$$

where  $\tau$  is experimentally set to 0.07. Consider Fig. 2 and assume, for example, that detection token 4 is selected to match the groundtruth (red) bounding box by the Hungarian matching algorithm. In this case,  $L^{OTA}$  encourages the embeddings  $\hat{o}_0^l$  and  $\hat{o}_1^l$  of the words “fuzzy” and “bench” to be closer to the embedding of  $\hat{o}_4^{det}$  than those of the words “closest,” “to,” and “you”.

Similarly, let  $\mathcal{O}_i \subseteq \{j : 1 \leq j \leq g\}$  be the set of object indices that a projected language token feature  $\hat{o}_i^l$  should be aligned to. This alignment is encouraged by the text-object alignment loss  $L^{TOA}$

$$L^{TOA} = \sum_{i=1}^m \frac{1}{|\mathcal{O}_i|} \sum_{j \in \mathcal{O}_i} -\log \left( \frac{\exp(\hat{o}_i^l \hat{o}_j^{det} / \tau)}{\sum_{k=1}^g \exp(\hat{o}_i^l \hat{o}_k^{det} / \tau)} \right), \quad (8)$$

which is a symmetric version of  $L^{OTA}$ . The overall object alignment loss is  $L^{OA} = \frac{1}{2}L^{OTA} + \frac{1}{2}L^{TOA}$ .

**Patch-Text Alignment Loss.** While the object-text alignment loss is desirable to encourage feature space alignment between *predicted* boxes and corresponding text tokens, this assumes that the predicted boxes are correct. However, as illustrated in Fig. 3, this assumption is frequently violated during training, especially in the early epochs. In this case, the object-text alignment loss can provide ambiguous supervision (e.g. supervision from background instead of the object in the early epochs of Fig. 3). To address this limitation, we propose to align the features extracted from the *ground truth patches at the input of the network* (green patches of Fig. 2) with the features of the associated text tokens. More specifically, the indices of input patches ( $p_5$  and  $p_6$  in Fig. 2) that have at least 0.5 IOU with the ground truth box (red box in Fig. 2) are first identified and the corresponding features aligned with those of the ground truth language tokens (“fuzzy” and “bench” in Fig. 2). As shown in Fig. 3, while each patch feature may not provide full information about the object, the patch features never provide contradictory supervision. This helps stabilize the learning.

To align text and patch features, the text transformer output  $o^l$  and vision transformer outputs  $o^v$  are first projected to a shared feature space using the

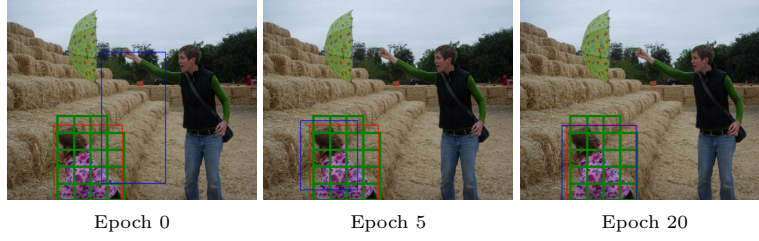


Fig. 3: **Visualization of predicted box** through training epochs. Red, green and blue boxes indicate ground truth box, ground truth patch and predicted box respectively. Patch-Text Alignment loss helps localization. Best viewed in color digitally.

		$p_1$	$p_2$	$p_3$	$p_4$	$p_5$	$p_6$
fuzzy	$t_1$	0	0	0	0	1	1
bench	$t_2$	0	0	0	0	1	1
closest	$t_3$	0	0	0	0	0	0
to	$t_4$	0	0	0	0	0	0
you	$t_5$	0	0	0	0	0	0

(a)

		$p_1$	$p_2$	$p_3$	$p_4$	$p_5$	$p_6$
fuzzy	$t_1$	0	0	0	0	0.5	0.5
bench	$t_2$	0	0	0	0	0.5	0.5
closest	$t_3$	0	0	0	0	0	0
to	$t_4$	0	0	0	0	0	0
you	$t_5$	0	0	0	0	0	0

(b)

		$p_1$	$p_2$	$p_3$	$p_4$	$p_5$	$p_6$
fuzzy	$t_1$	0	0	0	0	0.5	0.5
bench	$t_2$	0	0	0	0	0.5	0.5
closest	$t_3$	0	0	0	0	0	0
to	$t_4$	0	0	0	0	0	0
you	$t_5$	0	0	0	0	0	0

(c)

Fig. 4: **Patch-Text Alignment Loss**: (a) Ground truth for patch-text alignment. (b) Column wise normalization (c) Row wise normalization.

patch alignment projection heads  $\mathcal{H}_l^{PA}$  and  $\mathcal{H}_v^{PA}$  respectively, which leads to  $\tilde{o}^l$  and  $\tilde{o}^v$ . Given  $m$  text tokens and  $n$  visual patches, a ground truth table  $A \in \mathbb{R}^{m \times n}$  is then created, where  $A_{ij}$  is 1 if text token  $i$  and visual patch  $j$  should be aligned and 0 otherwise, as shown in Fig. 4(a). This matrix is column normalized into matrix  $A^C$  (Fig. 4 (b)) using

$$A_{ij}^c = \begin{cases} \frac{A_{ij}}{\sum_i A_{ij}} & \text{if } i \in \mathcal{S}^{TPA} \\ 0 & \text{otherwise} \end{cases}, \quad (9)$$

where  $\mathcal{S}^{TPA} = \{i | \sum_j A_{ij} > 0\}$  is the set of token indices associated with at least 1 patch. For example,  $\mathcal{S}^{TPA} = \{1, 2\}$  for Fig. 4. The row-wise normalization of  $A$  into matrix  $A^r$  is defined in a similar fashion, as shown in Fig. 4 (c).

To encourage the fine-grained alignment between patches and text, a text-patch alignment loss  $L^{TPA}$  is defined as

$$L^{TPA} = \sum_i \sum_j p_{ij}^c \ln \frac{p_{ij}^c}{A_{ij}^c} \quad \forall i \in \mathcal{S}^{TPA}, \quad (10)$$

where  $p_{ij}^c = \frac{\exp(\tilde{o}_i^{l\top} \tilde{o}_j^v) / \tau}{\sum_k \exp(\tilde{o}_i^{l\top} \tilde{o}_k^v) / \tau}$  is a contrastive distribution defined over the language and visual embedding. For example, the word “fuzzy” ( $t_1$ ) of Fig. 4(b) should be aligned with patch  $p_5$  and  $p_6$  according to the ground truth distribution  $A_1^c = [0; 0; 0; 0; 0.5; 0.5]$ . (10) minimizes the KL divergence between this and the predicted distribution  $p_1^c = [p_{11}^c; \dots; p_{1n}^c]$ . This encourages the alignment of the feature vector  $\tilde{o}_1^l$  extracted from word  $t_1$  with the feature vectors  $\tilde{o}_5^v$  and  $\tilde{o}_6^v$  extracted from patches  $p_5$  and  $p_6$ .

Table 1: **Architecture Design Ablation:** Comparison between YORO and ablated versions missing either learnable *det* tokens (w/o det) or *cls* feature  $o_{cls}^l$  (w/o cls).

Model Variant	Refcoco			Refcoco+			Refcocog	
	Val	TestA	TestB	Val	TestA	TestB	Val	Test
YORO w/o det.	74.9 (-8.0)	79.4 (-6.2)	68.8 (-8.6)	66.6 (-6.9)	73.2 (-5.4)	59.2 (-5.7)	69.7 (-3.7)	68.9 (-5.4)
YORO w/o cls.	82.4 (-0.5)	84.8 (-0.8)	76.7 (-0.7)	72.8 (-0.7)	77.7 (-0.9)	64.0 (-0.9)	72.4 (-1.0)	73.9 (-0.4)
YORO	<b>82.9</b>	<b>85.6</b>	<b>77.4</b>	<b>73.5</b>	<b>78.6</b>	<b>64.9</b>	<b>73.4</b>	<b>74.3</b>

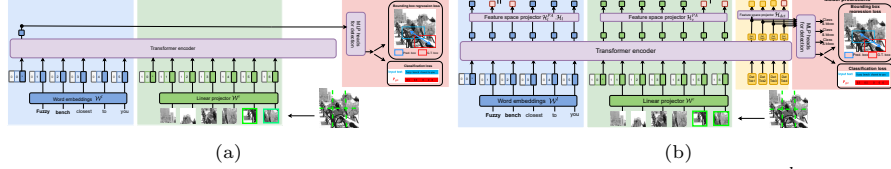


Fig. 5: (a) YORO w/o det. token. (b) YORO w/o cls. feature  $o_{cls}^l$ .

Similar to the  $L^{TPA}$  of (10), the patch-text alignment loss  $L^{PTA}$  aligns each patch to its associated tokens. As shown in Fig. 4(c),  $p_6$  should be aligned to  $t_1$  and  $t_2$ . The final patch alignment loss  $L^{PA}$  is  $L^{PA} = \frac{1}{2}L^{TPA} + \frac{1}{2}L^{PTA}$ .

**Combined Loss** The final loss combines the bounding box regression loss, the classification loss, the object alignment loss and the patch alignment loss with  $L = L^{bbox} + L^{cls} + L^{OA} + L^{PA}$ .

## 5 Experiments

### 5.1 Dataset

YORO is evaluated on five VG datasets using the accuracy metric, where the predicted box is correct when the IOU with the ground truth box is at least 0.5. **RefCoco/RefCoco+/RefCocog** [30, 80] are curated from MSCOCO [39]. RefCOCO consists of 142,209 referred expressions for 50,000 objects in 19,994 images. We adopt the split of [8, 85, 72, 78, 41, 5, 67] into train, val, testA, and testB datasets with 120,624, 10,834, 5,657 and 5,095 expressions respectively. RefCOCO+ contains 141,564 expressions for 49,856 objects in 19,992 images. We use the split of [88, 40, 4, 48, 78, 18], where train, val, testA and testB datasets contain 120,191, 10,758, 5,726 and 4,889 image text pairs, respectively. RefCOCog consists of 85,474 image text pairs of 54,822 objects across 26,711 images. We adopt the UMD split of [79, 43, 40, 75, 48], into 80,512, 4,896 and 9,602 pairs for train, val and test dataset, respectively.

**ReferItGame** [30] contains 19,997 images from the SAIAPR-12 [15] dataset. It has 130,363 text expressions for 99,296 objects. We follow the split of [24, 85, 49, 79, 71] for training, testing and validation set.

**CopsRef** [6], derived from GQA [27], contains 508 object classes and average expression length 14.4. When compared to RefCOCO and RefCOCog (80 object classes and average length around 6), CopsRef has longer expression and more diverse object classes. We adopt the **WithoutDist** version of the dataset [6, 14], where train, val and test contains 119,628, 12,586 and 16,524 pairs respectively.

Table 2: Loss Ablation. CL: Classification loss, RE: Bounding box regression loss, OA: Object-text Alignment Loss and PA: Patch-text Alignment loss.

Loss				CopsRef		ReferItGame	
CL	RE	OA	PA	Val.	Test	Val.	Test
✓	✓			67.05	70.78	71.55	69.86
✓	✓	✓		67.73 (+0.68)	71.0 (+0.22)	71.89 (+0.34)	70.39 (+0.53)
✓	✓	✓	✓	<b>68.08</b> (+1.03)	<b>71.3</b> (+0.52)	<b>72.67</b> (+1.12)	<b>71.9</b> (+2.04)

Table 3: Ablation of number of detection tokens on RefCoco+.

Det.	Tok. #	RefCoco+		
		Val.	TestA	TestB
1		73.2	<b>79.1</b>	<b>64.9</b>
5		<b>73.5</b>	78.6	<b>64.9</b>
10		72.8	78.78	62.8

## 5.2 Training and Evaluation Details

We summarize some implementation details of YORO here to enable reproducibility<sup>3</sup>. YORO leverages the **bert-base-uncased** [10] pre-trained tokenizer from HuggingFace [73] with a maximum text length of 40. The pre-trained checkpoint from [32] is used to obtain initial weights. Five detection tokens are used (see Table 3 for a token number ablation) and randomly initialized using a normal distribution of zero mean and standard deviation 0.02. The AdamW optimizer [45] with initial learning rate of  $10^{-4}$  and weight decay of  $10^{-2}$  is adopted. This warm up setting is used for the first 10% training epochs and the learning rate is then linearly decayed to 0. YORO is pre-trained on the concatenated detection dataset curated by [29] for 40 epochs. Since the VG task contains a single referred object, image-text pairs with a single bounding box are sampled. This pre-training allows a good initialization of the detection branch. The pre-trained model is fine-tuned on downstream datasets for 40 epochs. All experiments are conducted using PyTorch [51] with batch size 128.

## 5.3 Ablation study

**Architecture Design.** The performance of YORO is compared to the two variants of Fig. 5 on RefCoco/RefCoco+/RefCocog, with the results of Table 1.

**YORO w/o det. token.** Fig. 5(a) shows a simple extension of the ViLT [32] architecture to VG. In this case, the classification feature of the text modality  $o_{cls}^l$  (See Sec. 3.3) is forwarded through the detection head to regress a bounding box, without the use of detection branch or any detection tokens. The first row of Table 1 shows that, when compared to YORO, this is between 3.7% and 8.6% worse. For almost all splits, the performance loss is larger than 5%, demonstrating a clear benefit in using detection tokens for the VG task.

**YORO w/o cls. feature  $o_{cls}^l$ .** Conversely, the use of detection tokens without the holistic classification feature  $o_{cls}^l$  is implemented with the architecture of Fig. 5(b). This performs slightly worse (0.4% to 1%) than YORO. The holistic feature contributes to YORO’s performance but is much less critical than the detection tokens. Altogether, these ablations support the proposed design.

**Loss Function.** Table 2 studies the importance of the various losses in Sec. 4. The baseline consists of classification (CL) and bounding-box regression (RE) losses. Adding object-text (OA) alignment loss improves accuracy by +0.22% and +0.53%. When the PA loss is optimized, the gain over baseline increases to +0.52% and +2.04% on the test set of CopsRef and ReferItGame respectively.

<sup>3</sup> Code available at <https://github.com/chihhuiho/yoro>



Table 4: **Refcoco/Refcoco+/Refcocog**. Compared with single-stage models, YORO achieves the best accuracy. YORO is one of the smallest models and is also the fastest.

Method	Backbone		Refcoco			Refcoco+			Refcocog		Param. (M)	FPS
	Lang.	Visual	Val	TestA	TestB	Val	TestA	TestB	Val	Test		
FAOA [77]	Bert	Darknet53	72.05	74.81	67.91	-	-	-	-	-	182.9	18.9
RCCF [38]	BiLSTM	DLA34	-	81.06	71.85	-	70.35	56.32	-	-	-	40
SSG [4]	BiLSTM	Darknet53	-	76.51	67.50	-	62.14	49.27	58.80	-	-	-
MCN [48]	BiGRU	Darknet53	80.08	82.29	74.98	67.16	72.86	57.31	66.46	66.01	-	-
ReSC [76]	Bert	Darknet53	76.59	78.22	73.25	63.23	66.64	55.53	64.87	64.87	179.9	15.8
TransVG [9]	Bert	Res101	81.02	82.72	<b>78.35</b>	64.82	70.70	56.94	68.67	67.73	149.7	18.1
<b>YORO</b>	Bert	Linear	<b>82.9</b>	<b>85.6</b>	77.4	<b>73.5</b>	<b>78.6</b>	<b>64.9</b>	<b>73.4</b>	<b>74.3</b>	<b>114.3</b>	<b>41.9</b>

Table 5: Compared to single/two-stage methods, YORO achieves SOTA on **ReferItGame** [30].

Two-stage	Acc.	Single-stage(*)	Acc.
CMN [24]	28.33	RCCF* [38]	63.79
VC [85]	31.13	SSG* [4]	54.24
Luo [49]	31.85	ReSC-B* [76]	64.33
MAttNet [79]	29.04	ReSC-L* [76]	64.60
Sim. Net [71]	34.54	FAOA* [77]	59.30
CITE [52]	35.07	ZSGNet* [64]	58.63
PIRC [33]	59.13	TransVG* [9]	70.73
DDPN [83]	63.00	<b>YORO *</b>	<b>71.90</b>

Table 6: YORO achieves SOTA on **Cop-sRef** [6] w/o using g.t. box.

Method (* indicates single stage)	Proposal	Acc.
Grounder [63]	g.t.	75.7
Shuffle-Grounder	g.t.	58.5
Obj-Attr-Grounder	g.t.	68.8
MattNet [79]	g.t.	77.9
CM-Att-Erase [43]	g.t.	80.4
MattNet-Mine [6]	g.t.	78.4
VGTR-ResNet50 [14]	Pred.	66.73
VGTR-ResNet101 [14]	Pred.	67.75
ReSC-B* [76]	Pred.	64.49
ReSC-L* [76]	Pred.	65.32
<b>YORO *</b>	Pred.	<b>71.3</b>

This indicates both OA and PA loss improve detection accuracy on multiple datasets. The improved performance of PA validates the importance of accounting for the inaccuracy of bounding box predictions in the early training epochs.

**Number of Detection Tokens.** Table 3 presents an ablation of the number of detection tokens. On average, using 1 or 5 tokens has comparable results, with a slight decrease for 10 tokens. Since the VG task only contains a single referent, the lack of benefit in using a large number of detection tokens is sensible.

#### 5.4 Quantitative Comparison

**RefCoco/RefCoco+/RefCocog.** Table 4 summarizes performance of single-stage methods on RefCoco, RefCoco+ and RefCocog for different splits (val, testA and testB). These are the methods that achieve real time speed (FPS > 15) and thus directly comparable to YORO. YORO outperforms all single stage methods on seven splits (between abs +1.88% and +8.6%). The only exception is Refcoco testB (-0.95% lower than TransVG [9]). The most competitive method in this class is TransVG, which has  $1.3\times$  the size, is  $2.3\times$  slower than YORO, and achieves significantly lower accuracies on most splits. The only method of speed comparable to YORO (RCCF) has significantly lower accuracy on all splits (up to a 14 point drop on Refcocog Val). A more extensive comparison, including the much heavier two stage models, is presented in the appendix. Overall, as summarized in Fig. 1(b) and Fig. 6, YORO has a significantly better trade-off between speed, parameter size and accuracy than all methods in the literature.

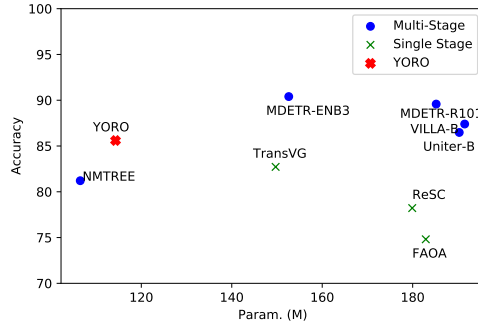


Fig. 6: YORO beats the prior Pareto front for the memory cost/accuracy trade-off.

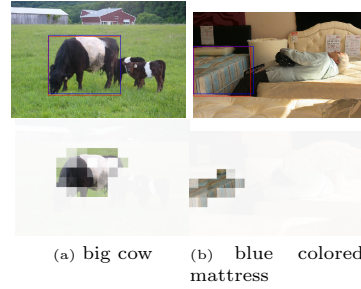


Fig. 7: Top: Predicted (blue) and ground truth (red) box. Bottom: Attention from the selected head.

**ReferItGame.** Table 5 compares YORO to both two-stage (left) and single-stage (right) baselines. YORO outperforms the best two-stage (DDPN [83]) and single-stage (TransVG [9]) method by 8.9% and 1.2% respectively. For roughly the same inference speed, YORO also beats RCCF [38] by 8.1%. Note that YORO achieves SOTA on ReferItGame for both speed and performance.

**CopsRef.** Table 6 summarizes the accuracy of YORO and the baselines on CopsRef [6]. The lower part of the table contains baselines that either detect the referent or select the referent from object proposals. YORO outperforms the previous SOTA by 3.4%, without using proposals. The upper part of the table refers to models that use *ground truth* bounding boxes. YORO fares well even under this unfair comparison, outperforming some of these baselines. Overall, these results demonstrate that YORO can generalize to VG datasets containing more object classes and longer input text.

### 5.5 Comparison of Trade-offs Between Size, Speed, and Accuracy

To quantify the efficiency of YORO, we compare YORO with other VG baselines in terms of parameter size and inference speed.

**Inference speed** (FPS) is measured by forwarding 100 images (batch size 1) from the Refcoco testA dataset through the VG model. The results are shown in the rightmost column of Table 4 and Fig. 1(b). All FPS measurements besides RCCF [38], MattNet [79] and DGA [75] were conducted by ourselves, using a single Titan Xp GPU and Intel Xeon CPU E5-2630 v4 @ 2.20G<sup>4</sup>. The FPS of [38, 79, 75] are copied from [38], which measures speed on a Titan Xp GPU (identical to ours) and Intel Xeon CPU E5-2680v4 @ 2.4G (superior to ours). YORO achieves 41.9 FPS and is the fastest VG model in the literature. It is between 3.3× and 14× faster than multi-stage models and at least twice as fast than all single-stage models other than RCCF, which has much weaker accuracy.

Fig. 1(b) summarizes the trade-off between speed and accuracy of several methods on Refcoco testA. Single-stage methods tend to have significantly lower accuracy than multi-stage methods. The major exception is YORO, which is

<sup>4</sup> Note that [77] reported a FPS of 26.3 for FAOA, using an NVIDIA 1080TI and Intel Core i9-9900K @ 3.60G set-up, while we measured a FPS of 18.9 under our set-up.

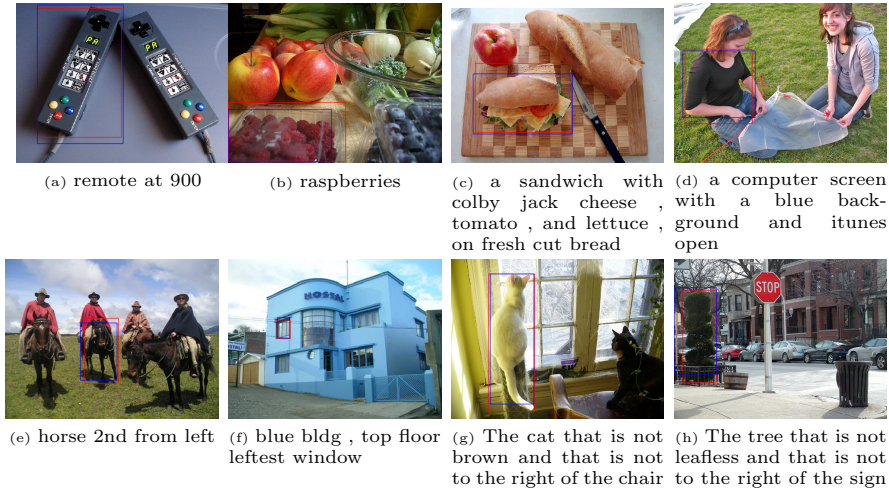


Fig. 8: **YORO Qualitative**: Visualization of the predicted box (blue) and ground truth box (red) in Refcoco+ (a-b), Refcocog (c-d), ReferItGame (e-f) and CopsRef (g-h). The input text for the referent is listed below each image.

substantially more accurate than the other single-stage methods, bridging a significant portion of the accuracy gap between the two types of models. Note that YORO is well above the line connecting MDETR and RCCF, which determines the Pareto front for the VG problem. It also has a large gain in inference speed (close to  $3.5\times$  faster) over MDETR. More plots are shown in the supplemental.

We also breakdown the percentage of inference time consumed per YOLO module. Multi-modal VG inputs (Sec. 3.1), transformer encoder (Sec. 3.2) and detection head (Sec. 3.4) consume 33.5%, 62.8% and 3.7% of the inference time, respectively. This suggests that future research should focus on optimizing the design of the transformer layers for both speed and accuracy.

**Parameter sizes** are compared in the second column from the right of Table 4 and Fig. 6. YORO is smaller than all baselines besides NMTREE [40], which uses a LSTM based language model. Note that, in particular, YORO is the smallest of the single-stage methods (See green  $\times$  in Fig. 6).

## 5.6 Qualitative Results

Fig. 7 shows the patches of higher attention weight for two example images, computed as follows. Let the detection token corresponding to the predicted box be token  $k$ . The attention weight between each patch and detection token  $k$  is converted to a heatmap and overlaid on the input image. It is clear that YORO attends to the referred objects. See appendix for other examples. Fig. 8 and the appendix show predictions by YORO on all 5 datasets. YORO can perform high quality detection not only on short sentences (a-b), but also on convoluted sentences (c-h). It also performs well in the presence of challenging input text, such as the digits of (a,e), abbreviation of (f), and objects of various sizes (b,f,h).

## 6 Conclusion

In this work, we investigate the speed accuracy trade-off of VG models. To pursue a better trade-off, we proposed YORO, a single-stage end-to-end trainable architecture. A novel patch-text alignment loss is proposed to improve the alignment between language and image features. Experiments show that YORO outperforms existing single-stage methods. It is the fastest VG and achieves the best trade-off between size, speed, and accuracy in the literature. We believe this is of interest for many resource constrained VG applications.

## References

1. Bochkovskiy, A., Wang, C.Y., Liao, H.Y.M.: YOLOv4: Optimal speed and accuracy of object detection. arXiv preprint arXiv:2004.10934 (2020) [2](#)
2. Carion, N., Massa, F., Synnaeve, G., Usunier, N., Kirillov, A., Zagoruyko, S.: End-to-end object detection with transformers. ArXiv **abs/2005.12872** (2020) [3](#), [5](#), [6](#), [26](#)
3. Chen, T., Deng, J., Luo, J.: Adaptive offline quintuplet loss for image-text matching. arXiv preprint arXiv:2003.03669 (2020) [8](#)
4. Chen, X., Ma, L., Chen, J., Jie, Z., Liu, W., Luo, J.: Real-time referring expression comprehension by single-stage grounding network. ArXiv **abs/1812.03426** (2018) [2](#), [4](#), [10](#), [12](#), [27](#)
5. Chen, Y.C., Li, L., Yu, L., Kholy, A.E., Ahmed, F., Gan, Z., Cheng, Y., Liu, J.: Uniter: Universal image-text representation learning. In: ECCV (2020) [2](#), [3](#), [8](#), [10](#), [27](#)
6. Chen, Z., Wang, P., Ma, L., Wong, K.Y.K., Wu, Q.: Cops-ref: A new dataset and task on compositional referring expression comprehension. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (June 2020) [10](#), [12](#), [13](#)
7. Cho, J., Lei, J., Tan, H., Bansal, M.: Unifying vision-and-language tasks via text generation. In: ICML (2021) [3](#), [4](#), [27](#)
8. Deng, C., Wu, Q., Wu, Q., Hu, F., Lyu, F., Tan, M.: Visual grounding via accumulated attention. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2018) [1](#), [10](#), [27](#)
9. Deng, J., Yang, Z., Chen, T., Zhou, W., Li, H.: Transvg: End-to-end visual grounding with transformers. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 1769–1779 (October 2021) [2](#), [4](#), [12](#), [13](#), [26](#), [27](#)
10. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. In: NAACL (2019) [5](#), [11](#)
11. Diao, H., Zhang, Y., Ma, L., Lu, H.: Similarity reasoning and filtration for image-text matching. In: AAAI (2021) [8](#)
12. Ding, H., Liu, C., Wang, S., Jiang, X.: Vision-language transformer and query generation for referring segmentation. In: Proceedings of the IEEE International Conference on Computer Vision (2021) [27](#)
13. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N.: An image is worth 16x16 words: Transformers for image recognition at scale. In: International Conference on Learning Representations (2021), <https://openreview.net/forum?id=YicbFdNTTy> [2](#), [4](#), [5](#), [6](#)
14. Du, Y., Fu, Z., Liu, Q., Wang, Y.: Visual grounding with transformers. CoRR **abs/2105.04281** (2021), <https://arxiv.org/abs/2105.04281> [3](#), [10](#), [12](#), [27](#)
15. Escalante, H.J., Hernández, C.A., Gonzalez, J.A., López-López, A., Montes, M., Morales, E.F., Enrique Sucar, L., Villaseñor, L., Grubinger, M.: The segmented and annotated iapr tc-12 benchmark. Comput. Vis. Image Underst. **114**(4), 419–428 (Apr 2010). <https://doi.org/10.1016/j.cviu.2009.03.008>, <https://doi.org/10.1016/j.cviu.2009.03.008> [10](#)
16. Fang, Y., Liao, B., Wang, X., Fang, J., Qi, J., Wu, R., Niu, J., Liu, W.: You only look at one sequence: Rethinking transformer in vision through object detection. arXiv preprint arXiv:2106.00666 (2021) [3](#), [5](#), [26](#)

17. Fukui, A., Park, D.H., Yang, D., Rohrbach, A., Darrell, T., Rohrbach, M.: Multi-modal compact bilinear pooling for visual question answering and visual grounding. In: *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*. pp. 457–468. Association for Computational Linguistics, Austin, Texas (Nov 2016). <https://doi.org/10.18653/v1/D16-1044>, <https://aclanthology.org/D16-1044> 3
18. Gan, Z., Chen, Y.C., Li, L., Zhu, C., Cheng, Y., Liu, J.: Large-scale adversarial training for vision-and-language representation learning. In: *NeurIPS (2020)* 2, 3, 10, 27
19. Goyal, Y., Khot, T., Summers-Stay, D., Batra, D., Parikh, D.: Making the v in vqa matter: Elevating the role of image understanding in visual question answering. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* pp. 6325–6334 (2017) 2
20. Han, K., Wang, Y., Chen, H., Chen, X., Guo, J., Liu, Z., Tang, Y., Xiao, A., Xu, C., Xu, Y., Yang, Z., Zhang, Y., Tao, D.: A survey on visual transformer. *ArXiv abs/2012.12556* (2020) 5
21. Hong, R., Liu, D., Mo, X., He, X., Zhang, H.: Learning to compose and reason with language tree structures for visual grounding. *IEEE transactions on pattern analysis and machine intelligence* (2019) 27
22. Howard, A.G., Sandler, M., Chu, G., Chen, L.C., Chen, B., Tan, M., Wang, W., Zhu, Y., Pang, R., Vasudevan, V., Le, Q.V., Adam, H.: Searching for mobilenetv3. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)* pp. 1314–1324 (2019) 2
23. Howard, A.G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., Andreetto, M., Adam, H.: Mobilenets: Efficient convolutional neural networks for mobile vision applications. *CoRR abs/1704.04861* (2017), <http://arxiv.org/abs/1704.04861> 1, 2
24. Hu, R., Rohrbach, M., Andreas, J., Darrell, T., Saenko, K.: Modeling relationships in referential expressions with compositional modular networks. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* pp. 4418–4427 (2017) 3, 10, 12, 27
25. Hu, R., Xu, H., Rohrbach, M., Feng, J., Saenko, K., Darrell, T.: Natural language object retrieval. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2016) 3
26. Huang, J., Rathod, V., Sun, C., Zhu, M., Balan, A.K., Fathi, A., Fischer, I.S., Wojna, Z., Song, Y., Guadarrama, S., Murphy, K.P.: Speed/accuracy trade-offs for modern convolutional object detectors. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* pp. 3296–3297 (2017) 2
27. Hudson, D.A., Manning, C.D.: Gqa: A new dataset for real-world visual reasoning and compositional question answering. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* pp. 6693–6702 (2019) 10
28. Jiao, L., Zhang, F., Liu, F., Yang, S., Li, L., Feng, Z., Qu, R.: A survey of deep learning-based object detection. *IEEE Access* 7, 128837–128868 (2019). <https://doi.org/10.1109/ACCESS.2019.2939201> 4
29. Kamath, A., Singh, M., LeCun, Y., Misra, I., Synnaeve, G., Carion, N.: Mdetr—modulated detection for end-to-end multi-modal understanding. *arXiv preprint arXiv:2104.12763* (2021) 3, 4, 8, 11, 27, 28
30. Kazemzadeh, S., Ordonez, V., Matten, M.A., Berg, T.L.: Referitgame: Referring to objects in photographs of natural scenes. In: *EMNLP (2014)* 10, 12
31. Khan, S.H., Naseer, M., Hayat, M., Zamir, S.W., Khan, F.S., Shah, M.: Transformers in vision: A survey. *ArXiv abs/2101.01169* (2021) 5

32. Kim, W., Son, B., Kim, I.: Vilt: Vision-and-language transformer without convolution or region supervision. In: Meila, M., Zhang, T. (eds.) Proceedings of the 38th International Conference on Machine Learning. Proceedings of Machine Learning Research, vol. 139, pp. 5583–5594. PMLR (18–24 Jul 2021), <http://proceedings.mlr.press/v139/kim21k.html> 2, 5, 6, 11
33. Kovvuri, R., Nevatia, R.: PIRC net : Using proposal indexing, relationships and context for phrase grounding. CoRR **abs/1812.03213** (2018), <http://arxiv.org/abs/1812.03213> 12
34. Kuhn, H.W.: The Hungarian Method for the Assignment Problem. Naval Research Logistics Quarterly **2**(1–2), 83–97 (March 1955). <https://doi.org/10.1002/nav.3800020109> 7
35. Lee, K.H., Chen, X., Hua, G., Hu, H., He, X.: Stacked cross attention for image-text matching. arXiv preprint arXiv:1803.08024 (2018) 8
36. Lewis, M., Liu, Y., Goyal, N., Ghazvininejad, M., Mohamed, A., Levy, O., Stoyanov, V., Zettlemoyer, L.: BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics. pp. 7871–7880. Association for Computational Linguistics, Online (Jul 2020). <https://doi.org/10.18653/v1/2020.acl-main.703>, <https://aclanthology.org/2020.acl-main.703> 4
37. Li, K., Zhang, Y., Li, K., Li, Y., Fu, Y.: Visual semantic reasoning for image-text matching. In: ICCV (2019) 8
38. Liao, Y., Liu, S., Li, G., Wang, F., Chen, Y., Qian, C., Li, B.: A real-time cross-modality correlation filtering method for referring expression comprehension. 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) pp. 10877–10886 (2020) 2, 4, 12, 13, 27
39. Lin, T.Y., Maire, M., Belongie, S.J., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: ECCV (2014) 10
40. Liu, D., Zhang, H., Wu, F., Zha, Z.J.: Learning to assemble neural module tree networks for visual grounding. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) (October 2019) 1, 10, 14, 27
41. Liu, J., Wang, L., Yang, M.H.: Referring expression generation and comprehension via attributes. In: 2017 IEEE International Conference on Computer Vision (ICCV). pp. 4866–4874 (2017). <https://doi.org/10.1109/ICCV.2017.520> 10, 27
42. Liu, L., Ouyang, W., Wang, X., Fieguth, P.W., Chen, J., Liu, X., Pietikäinen, M.: Deep learning for generic object detection: A survey. International Journal of Computer Vision **128**, 261–318 (2019) 4
43. Liu, X., Wang, Z., Shao, J., Wang, X., Li, H.: Improving referring expression grounding with cross-modal attention-guided erasing. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019) 1, 2, 10, 12, 27
44. Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., Stoyanov, V.: Roberta: A robustly optimized bert pretraining approach. ArXiv **abs/1907.11692** (2019) 4, 26
45. Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. In: International Conference on Learning Representations (2019), <https://openreview.net/forum?id=Bkg6RiCqY7> 11
46. Lu, J., Batra, D., Parikh, D., Lee, S.: Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In: Advances in Neural Information Processing Systems. pp. 13–23 (2019) 3, 27



47. Lu, J., Goswami, V., Rohrbach, M., Parikh, D., Lee, S.: 12-in-1: Multi-task vision and language representation learning. In: The IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (June 2020) [3](#), [27](#)
48. Luo, G., Zhou, Y., Sun, X., Cao, L., Wu, C., Deng, C., Ji, R.: Multi-task collaborative network for joint referring expression comprehension and segmentation. 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) pp. 10031–10040 (2020) [2](#), [4](#), [10](#), [12](#), [27](#)
49. Luo, R., Shakhnarovich, G.: Comprehension-guided referring expressions. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (July 2017) [1](#), [10](#), [12](#), [27](#)
50. Nagaraja, V.K., Morariu, V.I., Davis, L.S.: Modeling context between objects for referring expression understanding. In: ECCV (2016) [4](#)
51. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Köpf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., Chintala, S.: Pytorch: An imperative style, high-performance deep learning library. ArXiv [abs/1912.01703](#) (2019) [11](#)
52. Plummer, B.A., Kordas, P., Kiapour, M.H., Zheng, S., Piramuthu, R., Lazebnik, S.: Conditional image-text embedding networks. In: The European Conference on Computer Vision (ECCV) (2018) [12](#)
53. Qi, W., Yan, Y., Gong, Y., Liu, D., Duan, N., Chen, J., Zhang, R., Zhou, M.: Prophetnet: Predicting future n-gram for sequence-to-sequence pre-training. In: Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Findings. pp. 2401–2410 (2020) [4](#)
54. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., Sutskever, I.: Learning transferable visual models from natural language supervision. CoRR [abs/2103.00020](#) (2021), <https://arxiv.org/abs/2103.00020> [5](#)
55. Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., Liu, P.J.: Exploring the limits of transfer learning with a unified text-to-text transformer. Journal of Machine Learning Research **21**(140), 1–67 (2020), <http://jmlr.org/papers/v21/20-074.html> [4](#), [26](#)
56. Raffel, C., Shazeer, N.M., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., Liu, P.J.: Exploring the limits of transfer learning with a unified text-to-text transformer. ArXiv [abs/1910.10683](#) (2020) [4](#)
57. Ranftl, R., Bochkovskiy, A., Koltun, V.: Vision transformers for dense prediction. ArXiv preprint (2021) [5](#)
58. Redmon, J., Divvala, S., Girshick, R., Farhadi, A.: You only look once: Unified, real-time object detection. In: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 779–788 (2016). <https://doi.org/10.1109/CVPR.2016.911>, [2](#)
59. Redmon, J., Farhadi, A.: Yolo9000: Better, faster, stronger. 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR) pp. 6517–6525 (2017) [2](#)
60. Redmon, J., Farhadi, A.: Yolov3: An incremental improvement. ArXiv [abs/1804.02767](#) (2018) [2](#), [4](#)
61. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. In: Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 1. p. 91–99. NIPS’15, MIT Press, Cambridge, MA, USA (2015) [2](#), [3](#)

62. Rezatofghi, S.H., Tsoi, N., Gwak, J., Sadeghian, A., Reid, I.D., Savarese, S.: Generalized intersection over union: A metric and a loss for bounding box regression. 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) pp. 658–666 (2019) [7](#)
63. Rohrbach, A., Rohrbach, M., Hu, R., Darrell, T., Schiele, B.: Grounding of textual phrases in images by reconstruction. ArXiv **abs/1511.03745** (2016) [3](#), [12](#)
64. Sadhu, A., Chen, K., Nevatia, R.: Zero-shot grounding of objects from natural language queries. In: The IEEE International Conference on Computer Vision (ICCV) (October 2019) [12](#)
65. Sandler, M., Howard, A.G., Zhu, M., Zhmoginov, A., Chen, L.C.: Mobilenetv2: Inverted residuals and linear bottlenecks. 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition pp. 4510–4520 (2018) [2](#)
66. Strudel, R., Garcia, R., Laptev, I., Schmid, C.: Segmenter: Transformer for semantic segmentation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 7262–7272 (October 2021) [5](#)
67. Su, W., Zhu, X., Cao, Y., Li, B., Lu, L., Wei, F., Dai, J.: Vl-bert: Pre-training of generic visual-linguistic representations. In: International Conference on Learning Representations (2020), <https://openreview.net/forum?id=SygXPaEYvH> [2](#), [10](#), [27](#)
68. Suhr, A., Zhou, S., Zhang, A., Zhang, I., Bai, H., Artzi, Y.: A corpus for reasoning about natural language grounded in photographs. In: Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics. pp. 6418–6428. Association for Computational Linguistics, Florence, Italy (Jul 2019). <https://doi.org/10.18653/v1/P19-1644>, <https://aclanthology.org/P19-1644> [2](#)
69. Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., Jegou, H.: Training data-efficient image transformers & distillation through attention. In: International Conference on Machine Learning. vol. 139, pp. 10347–10357 (July 2021) [5](#)
70. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L.u., Polosukhin, I.: Attention is all you need. In: Guyon, I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R. (eds.) Advances in Neural Information Processing Systems. vol. 30. Curran Associates, Inc. (2017), <https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf> [4](#), [5](#), [6](#)
71. Wang, L., Li, Y., Huang, J., Lazebnik, S.: Learning two-branch neural networks for image-text matching tasks. TPAMI **41**(2), 394–407 (2019) [10](#), [12](#)
72. Wang, P., Wu, Q., Cao, J., Shen, C., Gao, L., Hengel, A.v.d.: Neighbourhood watch: Referring expression comprehension via language-guided graph attention networks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019) [1](#), [10](#), [27](#)
73. Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M., Davison, J., Shleifer, S., von Platen, P., Ma, C., Jernite, Y., Plu, J., Xu, C., Scao, T.L., Gugger, S., Drame, M., Lhoest, Q., Rush, A.M.: Transformers: State-of-the-art natural language processing. In: Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations. pp. 38–45. Association for Computational Linguistics, Online (Oct 2020), <https://www.aclweb.org/anthology/2020.emnlp-demos>. [6](#) [11](#)
74. Yang, S., Li, G., Yu, Y.: Cross-modal relationship inference for grounding referring expressions. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019) [1](#)

75. Yang, S., Li, G., Yu, Y.: Dynamic graph attention for referring expression comprehension. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) (October 2019) [1](#), [10](#), [13](#), [27](#)
76. Yang, Z., Chen, T., Wang, L., Luo, J.: Improving one-stage visual grounding by recursive sub-query construction. In: ECCV (2020) [2](#), [4](#), [12](#), [27](#)
77. Yang, Z., Gong, B., Wang, L., Huang, W., Yu, D., Luo, J.: A fast and accurate one-stage approach to visual grounding. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) (October 2019) [1](#), [2](#), [4](#), [12](#), [13](#), [27](#)
78. Yu, F., Tang, J., Yin, W., Sun, Y., Tian, H., Wu, H., Wang, H.: Ernie-vil: Knowledge enhanced vision-language representations through scene graphs. Proceedings of the AAAI Conference on Artificial Intelligence **35**(4), 3208–3216 (May 2021), <https://ojs.aaai.org/index.php/AAAI/article/view/16431> [10](#), [27](#)
79. Yu, L., Lin, Z., Shen, X., Yang, J., Lu, X., Bansal, M., Berg, T.L.: Mattnet: Modular attention network for referring expression comprehension. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2018) [1](#), [2](#), [3](#), [10](#), [12](#), [13](#), [27](#)
80. Yu, L., Poirson, P., Yang, S., Berg, A.C., Berg, T.L.: Modeling context in referring expressions. ArXiv **abs/1608.00272** (2016) [10](#)
81. Yu, L., Tan, H., Bansal, M., Berg, T.L.: A joint speaker-listener-reinforcer model for referring expressions. 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR) pp. 3521–3529 (2017) [1](#)
82. Yu, L., Tan, H., Bansal, M., Berg, T.L.: A joint speaker-listener-reinforcer model for referring expressions. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 3521–3529 (2017). <https://doi.org/10.1109/CVPR.2017.375> [3](#), [27](#)
83. Yu, Z., Yu, J., Xiang, C., Zhao, Z., Tian, Q., Tao, D.: Rethinking diversified and discriminative proposal generation for visual grounding. ArXiv **abs/1805.03508** (2018) [1](#), [4](#), [12](#), [13](#), [27](#)
84. Zaheer, M., Guruganesh, G., Dubey, K.A., Ainslie, J., Alberti, C., Ontanon, S., Pham, P., Ravula, A., Wang, Q., Yang, L., et al.: Big bird: Transformers for longer sequences. Advances in Neural Information Processing Systems **33** (2020) [4](#)
85. Zhang, H., Niu, Y., Chang, S.F.: Grounding referring expressions in images by variational context. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2018) [1](#), [10](#), [12](#), [27](#)
86. Zhang, J., Zhao, Y., Saleh, M., Liu, P.J.: Pegasus: Pre-training with extracted gap-sentences for abstractive summarization (2019) [4](#)
87. Zhu, X., Su, W., Lu, L., Li, B., Wang, X., Dai, J.: Deformable {detr}: Deformable transformers for end-to-end object detection. In: International Conference on Learning Representations (2021), <https://openreview.net/forum?id=gZ9hCDWe6ke> [5](#)
88. Zhuang, B., Wu, Q., Shen, C., Reid, I., van den Hengel, A.: Parallel attention: A unified framework for visual object discovery through dialogs and queries. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2018) [1](#), [10](#), [27](#)

## Appendix - Supplementary material

### A Limitation and Broader Impact

In this work, we designed a visual grounding (VG) model which has less parameters than state-of-the-art and has inference speed much faster than real-time (15 fps). YORO is the fastest VG methods in the literature and out-performs single-stage VG approaches by large margins. It also achieves better speed/accuracy trade-off than multi-stage VG models with  $3.3\times$  to  $14\times$  faster speed.

We believe our approach being light-weight will facilitate visual grounding (VG) in devices with less computation resources, such as edge devices, mobile robots, low-powered devices (e.g. Raspberry-Pi), and AR/VR devices. In addition, the proposed model requires less computation resources which is likely to be environmental friendly. These aspects warrant further research and consideration.

### B Patch-Text Supervision

Figure 9 shows the predicted bounding box (blue) and ground truth box (red) of a same image across different epochs. It can be observed that the predicted box gradually becomes more accurate from epoch 0 to epoch 40 and shifts around various locations in different epochs. On the contrary, the green ground truth patches, that have  $\text{IOU} \geq 0.5$  with the ground truth box, stay consistent throughout the training epochs and make the proposed patch-text alignment loss a more stable. This shows the need for proposed novel patch-text loss.

### C Parameter Size / Accuracy and Speed /Accuracy Plots

More comparisons on Parameter size vs Accuracy and Inference Speed vs Accuracy are provided here (w.r.t. Sec. 5 in main paper). We hope the community is inspired to consider other factors which make the technology useful (like inference speed and parameters) and not just on only Accuracy and the race to beat state-of-the-art in performance metrics.

Figure 10, 11 and 12 illustrate the trade-off between the accuracy and the inference speed in terms of FPS for RefCoco, RefCoco+ and RefCocog, respectively. YORO achieves better trade-off between the speed and accuracy in different datasets.

Figure 13, 14 and 15 illustrate the trade-off between the accuracy and the number of parameter used by each baseline for RefCoco, RefCoco+ and RefCocog, respectively. YORO achieves competitive accuracy with smaller model size.

This indicates YORO is more suitable for edge devices used in robotic applications and AR/VR devices.

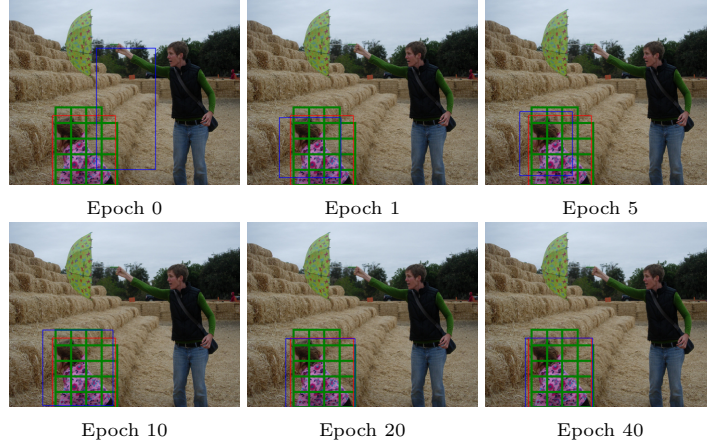


Fig. 9: **Patch-Text Supervision:** Need for consistent supervision for patch-alignment loss. Ground truth box is marked in red and green patches have overlapped with red box with  $\text{IOU} \geq 0.5$ . Blue box is the prediction across different epochs.

## D Qualitative Results

More qualitative visualizations from YORO are shown here (w.r.t Sec. 5 in main paper).

**Correct Examples.** Figure 16 contains correct detection results from YORO output. Examples from RefCoco and RefCoco+ have shorter query text, while examples from RefCocog are associated to more descriptive text. YORO can handle query text of various length and correctly detect the referent when the query contains digits. While the odd rows contain the predicted bounding box, the even rows show the patches that higher attention weight. More specifically, assuming the  $k^{\text{th}}$  detection token corresponds to the predicted bounding box, we extract the transformer attention weight between each patch and the  $k^{\text{th}}$  detection token. The attention weight is then converted to heatmap and overlaid on the input image. Take Figure 16(d) for example. Most patches around the bottom orange have higher attention weight.

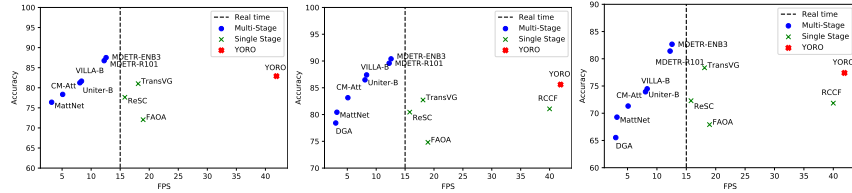


Fig. 10: Speed Accuracy Plot on RefCoco val, testA and testB dataset.

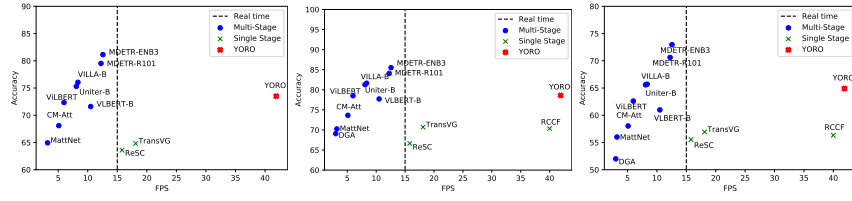


Fig. 11: Speed Accuracy Plot on RefCoco+ val, testA and testB dataset.

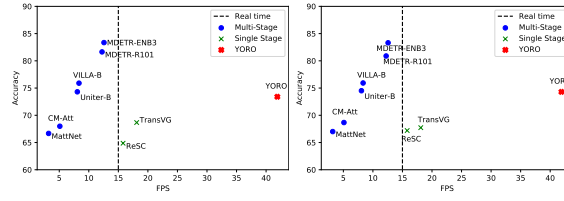


Fig. 12: Speed Accuracy Plot on RefCocog val, test dataset.

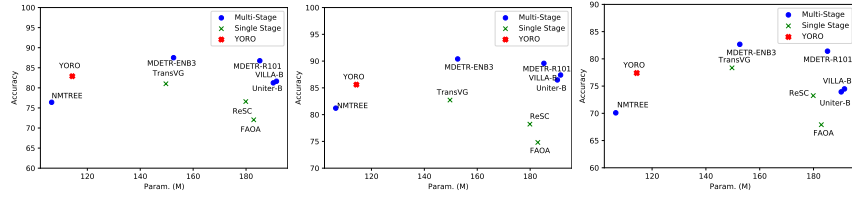


Fig. 13: Memory Accuracy Plot on RefCoco val, testA and testB dataset

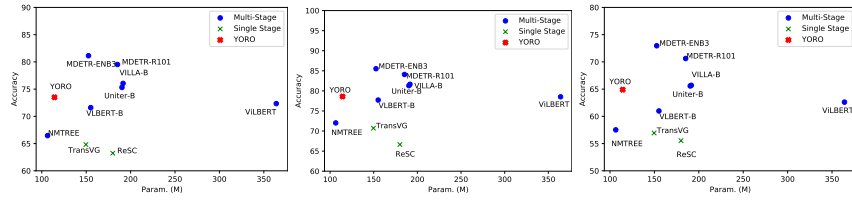


Fig. 14: Memory Accuracy Plot on RefCoco+ val, testA and testB dataset

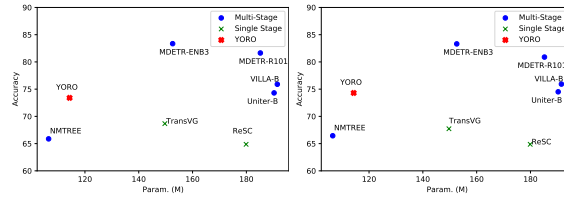


Fig. 15: Memory Accuracy Plot on RefCocog val, test dataset

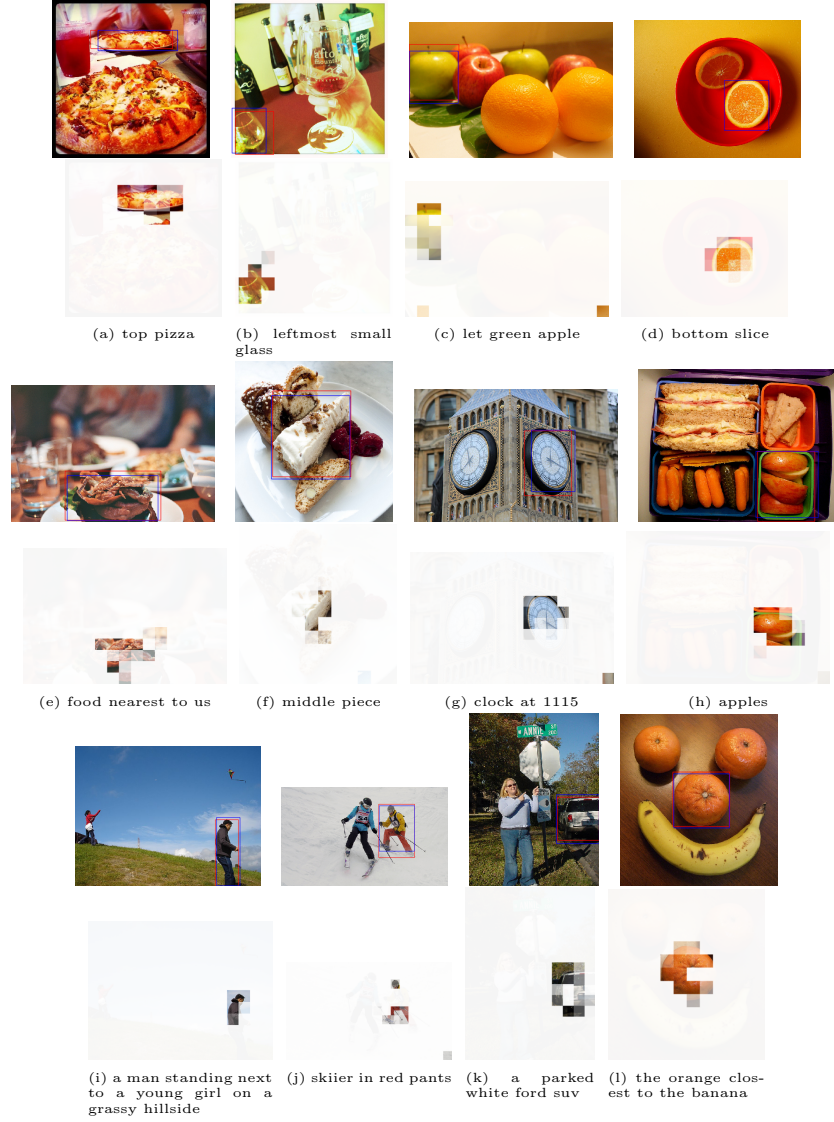


Fig. 16: Visualization of the predicted box (blue) and ground truth box (red) in Refcoco (a-d)/Refcoco+(e-h)/Refcocog(i-l) dataset. The patches that have higher attention are highlighted in the even rows.



## E Full Comparison with State of the Art

Table 7 extends the comparison of YORO with both single stage and multi-stage methods. Note that comparing YORO with multi-stage methods in terms of accuracy is **not an apples-to-apples comparison**, because multi-stage methods tends to sacrifice speed for accuracy improvement with the use of resource consuming models. The table is divided into blocks with the top performance within each block in bold.

**(Bi)LSTM language model Two-stage methods:** These methods are based on small (LSTM-based) language models of relatively few parameters that severely under-perform transformer-based approaches. Note that, for these models, parameter size does not correlate with inference speed, due to recursive LSTM computations. In fact, these are the slowest models considered. YORO significantly out-performs all these models in most splits (up +7 abs points), for marginally higher memory complexity (7% increase) and much faster inference (10× speed up).

**BERT language model Two-stage methods:** The main difference between YORO and these methods is the elimination of the visual backbone. This drastically reduces parameter size (between 60% and 74% of the parameters) and boosts inference speed by a substantial amount (between 4× and 7× speed-up). Despite these gains, YORO is comparable to the best of these methods, achieving superior performance on Refcoco val/testB (between +1.25 and +2.9) and comparable in other splits.

**Encoder-Decoder multi-stage methods:** Like YORO, these methods are transformer-based but use more powerful encoder-decoder models. While it has been argued that encoder-decoder models have more powerful predictive capability than encoder-only approaches [55, 2, 16], YORO outperforms all the methods in this class other than MDETR by large margins (between abs +3.0% and +10.1%). The comparison to MDETR is more complex, because it also uses the more powerful RoBERTa [44] language model. This enables MDETR to achieve higher accuracies, but also makes it a lot larger (1.3×) and slower (3.4×) than YORO.

**Single-stage methods:** These are the methods directly comparable to YORO. YORO outperforms all these methods on seven splits (between abs +1.88% and +8.6%). The only exception is Refcoco testB (-0.95% lower than TransVG [9]). The most competitive method in this class is TransVG, which has 1.3× the size and is 2.3× slower than YORO but achieves significantly lower accuracies on most of the splits. On the other hand, the only method of speed comparable to YORO (RCCF) has significantly lower accuracy on all splits (up to a 14 point drop on Refcocog Val). Overall, YORO has a significantly better trade-off between speed, parameter size and accuracy than all these methods.

Table 7: **Refcoco**, **Refcoco+** and **Refcocog**. YORO achieves SOTA performance when compared with single-stage models. YORO is one of the smallest models and is also the fastest. **Bold type** indicates the best model in each block.

Method	Backbone		Refcoco			Refcoco+			Refcocog		Param.	FPS
	Lang.	Visual	Val	TestA	TestB	Val	TestA	TestB	Val	Test	(M)	
<i><b>Two-Stage</b> Vision-Language methods based on (Bi)LSTM for Language Model</i>												
S.L.R. [82]	LSTM	VGG16	-	72.94	62.98	-	58.68	47.68	-	-	-	-
Luo [49]	BiLSTM	VGG16	-	67.94	55.18	-	57.05	43.33	-	-	-	-
Deng [8]	LSTM	VGG16	<b>81.27</b>	81.17	<b>80.01</b>	65.56	68.76	60.63	-	-	-	-
VC [85]	LSTM	VGG16	-	73.33	67.44	-	58.40	53.18	-	-	-	-
LGRAN [72]	LSTM	VGG16	-	76.6	66.4	-	64.0	53.4	-	-	-	-
Liu [41]	LSTM	VGG19	-	72.08	57.29	-	57.97	46.20	-	-	-	-
MattNet [79]	BiLSTM	Res101	76.40	80.43	69.28	64.93	70.26	56.00	66.67	67.01	-	3.2
CM-Att-Erase [43]	LSTM	Res101	78.35	<b>83.14</b>	71.32	<b>68.09</b>	<b>73.65</b>	<b>58.03</b>	<b>67.99</b>	<b>68.67</b>	-	5.1
CMN [24]	LSTM	VGG16	-	71.03	65.77	-	54.32	47.76	-	-	-	-
DDPN [83]	LSTM	Res101	76.8	80.1	72.4	64.8	70.5	54.1	-	-	-	-
DGA [75]	BiLSTM	VGG16	-	78.42	65.53	-	69.07	51.99	-	63.28	-	3
PLAN [88]	LSTM	VGG16	-	75.31	65.52	-	61.34	50.86	-	-	-	-
RvGTree [21]	BiLSTM	Res101	75.06	78.61	69.85	63.51	67.45	56.66	66.95	66.51	-	-
NMTREE [40]	BiLSTM	Res101	76.41	81.21	70.09	66.46	72.02	57.52	65.87	66.44	106.5	-
<i><b>Two-Stage</b> Vision-Language methods based on Bert-base for Language Model</i>												
Uniter [5]	Bert-B	Res101	81.24	86.48	73.94	75.31	81.3	65.58	74.31	74.51	190.3	8.1
VLBERT [67]	Bert-B	Res101	-	-	-	71.6	77.72	60.99	-	-	155.1	10.5
VILLA [18]	Bert-B	Res101	<b>81.65</b>	<b>87.4</b>	<b>74.48</b>	<b>76.05</b>	<b>81.65</b>	<b>65.7</b>	<b>75.9</b>	75.93	191.5	8.3
ERNIE ViLL [78]	Bert-B	Res101	-	-	-	74.02	80.33	64.74	-	-	-	-
12 in 1 [47]	Bert-B	ResXT152	-	80.58	-	-	73.25	-	-	<b>75.96</b>	-	-
ViLBERT [46]	Bert-B	Res101	-	-	-	72.34	78.52	62.61	-	-	364	6
<i><b>Multi-Stage</b> Vision-Language methods - Encoder-Decoder Transformers.</i>												
<b>*Note:</b> these approaches are more powerful than encoder-only models.												
VGTR [14]	Bi-LSTM	Res50	78.29	81.49	72.38	63.29	70.01	55.64	64.19	64.01	-	-
VGTR [14]	Bi-LSTM	Res101	79.20	82.32	73.78	63.91	70.09	56.51	65.73	67.23	-	-
VL-T5 [7]	T5	Res101	-	-	-	-	-	-	-	71.3	304.9	6.7
VLT [12]	RNN	Res50	76.20	80.31	71.44	64.19	68.40	55.84	61.03	60.24	-	-
MDETR [29]	Roberta-B	Res101	86.75	89.58	81.41	79.52	84.09	70.62	81.64	80.89	185.2	12.24
MDETR [29]	Roberta-B	ENB3	<b>87.51</b>	<b>90.4</b>	<b>82.67</b>	<b>81.13</b>	<b>85.52</b>	<b>72.96</b>	<b>83.35</b>	<b>83.31</b>	152.6	12.57
<i><b>Single-Stage</b> Vision-Language methods</i>												
FAOA [77]	Bert	Darknet53	72.05	74.81	67.91	-	-	-	-	-	182.9	18.9
RCCF [38]	BiLSTM	DLA34	-	81.06	71.85	-	70.35	56.32	-	-	-	40
SSG [4]	BiLSTM	Darknet53	-	76.51	67.50	-	62.14	49.27	58.80	-	-	-
MCN [48]	BiGRU	Darknet53	80.08	82.29	74.98	67.16	72.86	57.31	66.46	66.01	-	-
ReSC [76]	Bert-B	Darknet53	76.59	78.22	73.25	63.23	66.64	55.53	64.87	64.87	179.9	15.8
TransVG [9]	Bert-B	Res101	81.02	82.72	<b>78.35</b>	64.82	70.70	56.94	68.67	67.73	149.7	18.1
<b>YORO</b>	Bert-B	Linear	<b>82.9</b>	<b>85.6</b>	77.4	<b>73.5</b>	<b>78.6</b>	<b>64.9</b>	<b>73.4</b>	<b>74.3</b>	<b>114.3</b>	<b>41.9</b>

Pretraining	RefCoco			CopsRef	
	val	testA	testB	val	test
w/o	73.4	78	67.1	64.4	69.5
w/	82.9 (+9.5)	85.6 (+7.6)	77.4 (+10.3)	68.08 (+3.68)	71.3 (+1.8)

Table 8: Ablation w/ and w/o pre-training.

## F Ablation on Model Pretraining

As mentioned in the main paper, YORO is pre-trained on the concatenated detection dataset curated by [29] to allow a good initialization of the detection branch. Table 8 further highlights the importance of the pre-training on RefCoco and CopsRef dataset. The averaged gains are 9.13% and 2.74% on RefCoco and CopsRef respectively. This supports the need of pre-training stage for the detection branch.