# Variational Context: Exploiting Visual and Textual Context for Grounding Referring Expressions

Yulei Niu, Hanwang Zhang, Zhiwu Lu, Shih-Fu Chang

**Abstract**—We focus on grounding (i.e., localizing or linking) referring expressions in images, e.g., "largest elephant standing behind baby elephant". This is a general yet challenging vision-language task since it does not only require the localization of objects, but also the multimodal comprehension of context — visual attributes (e.g., "largest", "baby") and relationships (e.g., "behind") that help to distinguish the referent from other objects, especially those of the same category. Due to the exponential complexity involved in modeling the context associated with multiple image regions, existing work oversimplifies this task to pairwise region modeling by multiple instance learning. In this paper, we propose a variational Bayesian method, called Variational Context, to solve the problem of complex context modeling in referring expression grounding. Specifically, our framework exploits the *reciprocal relation* between the referent and context, *i.e.*, either of them influences estimation of the posterior distribution of the other, and thereby the search space of context can be greatly reduced. In addition to reciprocity, our framework considers the *semantic information* of context, *i.e.*, the referring expression can be reproduced based on the estimated context. We also extend the model to unsupervised setting where no annotation for the referent is available. Extensive experiments on various benchmarks show consistent improvement over state-of-the-art methods in both supervised and unsupervised settings.

**Index Terms**—Grounding referring expression, variational Bayesian model, referring expression generation

✦

## 1 INTRODUCTION

Grounding natural language in visual data is a hallmark of artificial intelligence, since it establishes a communication channel between humans, machines, and the physical world, underpinning a variety of multimodal artificial intelligence tasks such as robotic navigation [45], visual question answering [1], [21], [62], and visual chatbot [7]. Thanks to the rapid development in deep learning-based computer vision and natural language processing, we have witnessed promising results not only in grounding nouns (*e.g.*, object detection [37]), but also short phrases (*e.g.*, noun phrases [35] and relations [43], [59]). However, the more general task: grounding referring expressions [30], is still far from resolved due to the challenges in understanding of both language and scene compositions [13]. As illustrated in Fig. 1, given an input referring expression "largest elephant standing behind baby elephant" and an image with region proposals, a model that can only localize "elephant" is not satisfactory as there are multiple elephants. Therefore, the key for referring expression grounding is to comprehend the context. Here, we refer to *context* as the visual objects (*e.g.*, "elephant"), attributes (*e.g.*, "largest" and "baby"), and relationships (*e.g.*, "behind") mentioned in the expression that help to distinguish the referent from other objects.

One straightforward way of modeling the relations between the referent and context is to: 1) use external syntactic parsers to parse the expression into entities,
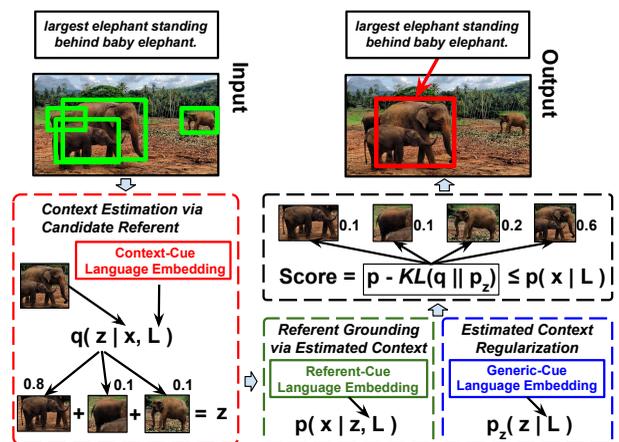


Fig. 1. The proposed Variational Context framework. Given an input referring expression and an image with region proposals, we localize the referent as output. We develop a grounding score function, with the variational lower-bound composed by three cue-specific multimodal modules, indicated by the description in the dashed color boxes.

modifiers, and relations [41], and then 2) apply visual relation detectors to localize them [59]. However, this two-stage approach is not practical due to the limited generalization ability of the detectors applied in the highly unrestricted language and scene compositions. To this end, recent approaches use multimodal embedding networks that jointly comprehend language and model the visual relations [15], [33]. Due to the prohibitively high cost of annotating both referent and context of referring expressions in images, multiple instance learning (MIL) [10] is usually adopted in them to handle the weak supervision of the unannotated context objects, by maximizing the joint likelihood of every region pair. However, for a referent, the

Y. Niu and Z. Lu are with the Beijing Key Laboratory of Big Data Management and Analysis Methods, School of Information, Renmin University of China, Beijing 100872, China. Email: niu@ruc.edu.cn, zhiwu.lu@gmail.com.
H. Zhang is with the School of Computer Science and Engineering, Nanyang Technological University, Singapore 639798. E-mail: hanwangzhang@gmail.com.
S.-F. Chang is with the Department of Electrical Engineering, Columbia University, New York, NY 10027, USA. Email: sfchang@ee.columbia.edu.

MIL framework essentially oversimplifies the number of context configurations of $N$ regions from $\mathcal{O}(2^N)$ to $\mathcal{O}(N)$. For example, to localize the "elephant" in Fig. 1, we may need to consider the other three elephants all together as a multinomial subset for modeling the context such as "largest", "behind" and "baby elephant".

In this paper, we propose a novel framework called *Variational Context* for grounding referring expressions in images. Compared to the previous MIL-based approaches [15], [33], our model approximates the combinatorial context configurations with weak supervision using a variational Bayesian framework [19]. Intuitively, it exploits the *reciprocity* between referent and context, given either of which can help to localize the other. As shown in Fig. 1, for each region $x$, we first estimate a coarse context $z$, which will help to refine the true localizations of the referent. This reciprocity is formulated into the variational lower-bound of the grounding likelihood $p(x|L)$, where $L$ is the text expression and the context is considered as a hidden variable $z$ (cf. Section 3). Specifically, the model consists of three multimodal modules: context posterior $q(z|x, L)$, referent posterior $p(x|z, L)$, and context prior $p_z(z|L)$, each of which performs a grounding task (cf. Section 4.3) that aligns image regions with a cue-specific language feature; each cue dynamically encodes different subsets of words in the expression $L$ that help the corresponding localization (cf. Section 4.2). In addition to reciprocity, our framework considers the *semantic information* of context, *i.e.*, the referring expression can be reproduced based on the referent and its estimated context. Specifically, the context prior $p_z(z|L)$ is resolved into the likelihood $p(L|z)$ for modeling the context-aware referring expression and the prior $p(z)$ following Bayes' theorem. Our proposed framework unifies both referring expression comprehension and generation to promote the context modeling.

Thanks to the reciprocity between referent and context, our model can not only be used in the conventional supervised setting, where there is annotation for referent, but also in the challenging unsupervised setting, where there is no instance-level annotation (*e.g.*, bounding boxes) of both referent and context. We perform extensive experiments on four benchmark referring expression datasets: RefCLEF [18], RefCOCO [57], RefCOCO+ [57], and RefCOCOg [30]. Our model consistently outperforms previous methods in both supervised and unsupervised settings. We also qualitatively show that our model can ground the context in the expression to the corresponding image regions (cf. Section 5). An earlier version of this work has appeared in [61]. The current paper 1) unifies referring expression comprehension and generation to promote complex context modeling, 2) updates state-of-the-art grounding results, 3) enriches qualitative results and failure cases studies, and 4) performs automatic evaluation and human evaluation on referring expression generation.

## 2 RELATED WORK

**Grounding Referring Expression**. Referring expression is the natural language description of a given object or region in an image. Grounding referring expression, also known as referring expression comprehension, intends to localize the target object given the referring expression. Different from grounding phrases [35], [36] and descriptive sentences [16], [39], the key for grounding referring expression is to use the context (or pragmatics in linguistics [44]) to distinguish the referent from other objects, usually of the same category [13]. However, most previous works resort to use holistic context such as the entire image [8], [16], [30], [39] or visual feature difference between regions [56], [57], [58]. Our model is similar to the works on explicitly modeling the referent and context region pairs [15], [33]. However, due to the lack of context annotation, they reduce the grounding task into a multiple instance learning framework [10]. As we will discuss later, this framework is not a proper approximation to the original task. There are also studies on visual relation detection that detect objects and their relationships [6], [22], [26], [56], [59], [60]. However, they are limited to a fixed-vocabulary set of relation triplets and hence are difficult to be applied in natural language grounding. Our cue-specific language feature is similar to the language modular network [15] that learns to decompose a sentence into referent/context-related words, which are different from other approaches that treat the expression as a whole [24], [28], [30], [58].

**Referring Expression Generation**. As the inverse task of grounding referring expression, referring expression generation [30] aims to produce a distinct expression about one object in an image. Different from image captioning [47], [53], referring expression generation mainly focuses on one specific object instead of the whole image. In addition, the generated referring expression should be unambiguous and include discriminative information attributes, location and relation. Most of early works have studied referring expression generation in neutral language processing [20], [31], [32], [46], [51]. However, they focused on the small and artificial datasets of past and have less concern about complex real-world vision problem. Recently, Kazemzadeh *et al.* [18] collected a large-scale dataset RefCLEF for natural pictures in real world via a two-player game, where one player provides a referring expression given the object in an image, and another player localizes the referent based on the referring expression and image. Other datasets RefCOCO, RefCOCO+ and RefCOCOg [30], [57] on MSCOCO were also collected in the same way. Based on the large-scale datasets, recent works make contributions to linking vision and language. The CNN-RNN structure, widely applied in image captioning, is generally used in referring expression generation. Recent works have investigated the combination of referring expression comprehension and generation. Our proposed Variational Context framework has the following differences in the combination strategy. 1) Compared to [30] that formulated expression generation as $\arg\max_L p(L|x, I)$ (where $L$ represents the description, $x$ represents the referent, and $I$ represents the image), we formulate the generation problem as $\arg\max_L p(L|x, z(x), I)$, by considering the context $z(x)$ of referent $x$. 2) Compared to [57] that proposed to jointly generate expressions of objects in the same class for unambiguous expression generation, we separately generate each expression of a referent with its context. 3) Compared to [58] that first accumulated comprehension and generation losses and then formulated the overall loss function as a multi-task

learning problem, our proposed framework first formulates visual grounding as a marginal distribution approximation problem, and then resolves the context prior $p(z|L)$ to include generation likelihood $p(L|z)$ in the lower-bound approximation.

**Variational Bayesian Model vs. Multiple Instance Learning**. Our proposed variational context framework is in a similar vein of the deep neural network based variational autoencoder (VAE) [19], which uses neural networks to approximate the posterior distribution of the hidden value $q(z|x)$, i.e., encoder, and the conditional distribution of the observation $p(x|z)$, i.e., decoder. VAE shows efficient and effective end-to-end optimization for the *intractable* log-sum likelihood $\log \sum_z p(x, z)$ that is widely used in generative processes such as image synthesis [55] and video frame prediction [54]. Considering the unannotated context as the hidden variable $z$, the referring expression grounding task can also be formulated into the above log-sum marginalization (cf. Eq. (2)). The MIL framework [10] is essentially a sum-log approximation of the log-sum, i.e., $\sum_z \log p(x, z)$. To see this, the max-pooling function $\log \max_z p(x, z)$ used in [15] can be viewed as the sum-log $\sum_z \log p(x|z)p(z)$, where $p(z) = 1$ if $z$ is the correct context and $0$ otherwise, indicating there is only one positive instance; maximizing the noisy-or function $\log(1 - \prod_z (1 - p(x, z)))$ used in [33] is equivalent to maximize $\sum_z \log p(x, z)$, assuming there is at least one positive instance. However, due to the numerical property of the log function, this sum-log approximation will unnecessarily force every $(x, z)$ pair to explain the data [11]. Instead, we use the variational Bayesian upper-bound to obtain a better sum-log approximation. Note that visual attention models [2], [53] simplify the variational lower bound by assuming $p(z) = q(z|x)$; however, we explicitly use the KL divergence $KL(q(z|x)||p(z))$ in the lower bound to regularize the approximate posterior $q(z|x)$ not being too far from the prior $p(z)$.

## 3 VARIATIONAL CONTEXT

In this section, we derive the variational Bayesian formulation of the proposed Variational Context framework and the objective function for training and test.

### 3.1 Problem Formulation

In this paper, we follow the classical definition of grounding referring expressions, where the region proposal generation stage is assumed to be done and only the region retrieval stage is considered. Note that one-stage visual grounding has the potential to be efficient [4], [9], which is out of the discussion in this paper. The classical task of grounding a referring expression $L$ in an image $I$ aims to retrieval the target object $x^*$ from a candidate set of regions $\mathcal{X}$. Formally, we maximize the log-likelihood of the conditional distribution to localize the referent region $x^* \in \mathcal{X}$:

$$x^* = \arg\max_{x \in \mathcal{X}} \log p(x|L), \qquad (1)$$

where we omit the image $I$ in $p(x|I, L)$.

As there is usually no annotation for the context, we consider it as a hidden variable $z$. Therefore, Eq. (1) can be rewritten as the following maximization of the log-likelihood of the conditional marginal distribution:

$$x^* = \arg\max_{x \in \mathcal{X}} \log \sum_z p(x, z|L). \qquad (2)$$

Note that $z$ is NOT necessary to be one region as assumed in recent MIL approaches [15], [33], i.e., $z \in \mathcal{X}$. For example, the contextual objects "surrounding elephants" in "a bigger elephant than the surrounding elephants" should be composed by a multinomial subset of $\mathcal{X}$, resulting in an extremely large sample space that requires $\mathcal{O}(2^{|\mathcal{X}|})$ search complexity. Therefore, the marginalization in Eq. (2) is intractable in general.

To this end, we use the variational lower-bound [19] to approximate the marginal distribution in Eq. (2) as:

$$\log p(x|L) = \log \sum_z p(x, z|L) \geq \mathcal{Q}(x, L) =$$

$$\underbrace{\mathbb{E}_{z \sim q_\phi(z|x,L)} \log p_\theta(x|z, L)}_{\text{Localization}} - \underbrace{KL(q_\phi(z|x, L)||p_\omega(z|L))}_{\text{Regularization}}, \quad (3)$$

where $KL(\cdot)$ is the Kullback-Leibler divergence, $\phi$, $\theta$, and $\omega$ are independent parameter sets for the respective distributions. As shown in Fig. 1, the lower bound $\mathcal{Q}(x, L)$ offers a new perspective for exploiting the reciprocal nature of referent and context in referring expression grounding.

#### 3.1.1 Localization

This term calculates the localization score for $x$ given an estimated context $z$, using the referent-cue of $L$ parameterized by $\theta$. In particular, we design a new posterior $q_\phi(z|x, L)$ that approximates the true context posterior $p(z|x, L)$, which models the context $z$ using the context-cue of $L$ parameterized by $\phi$. In the view of variational auto-encoder [19], [42], this term works in an encoding-decoding fashion: $q_\phi$ is the encoder from $x$ to $z$, and $p_\theta$ is the decoder from $z$ to $x$.

#### 3.1.2 Regularization

Since the Kullback-Leibler divergence ($KL$) is non-negative, maximizing $\mathcal{Q}(x, L)$ would encourage that the posterior $q_\phi$ is similar to the prior $p_\omega$, i.e., *the estimated context $z$ sampled from $q_\phi(z|x, L)$ should not be too far from the referring expression*, which is modeled by $p_\omega(z|L)$ with the generic-cue of $L$ parameterized by $\omega$. This term is necessary as the estimated context $z$ could be overfitted to region features that are inconsistent with the visual context described in the expression.

Furthermore, we hope that the estimated context $z$ contains necessary semantic information, which can helps to reproduce the referring expression. This can be realized by unifying both referring expression comprehension and generation, and resolving the context prior $p(z|L)$ as:

$$p(z|L) = g(x, L)p(L|z), \qquad (4)$$

where $g(x, L)$ is the function representing $p(z)/p(L)$, and we omit the referent region $x$ in $z(x)$ for simplicity. The likelihood $p(L|z)$ models the referring expression $L$ based on the estimated context $z$. Applying Eq. (4) to Eq. (3), we can get another lower bound $\mathcal{Q}'(x, L)$:

$$\mathcal{Q}'(x, L) = \mathbb{E}_{z \sim q_\phi(z|x,L)}[\log p_\theta(x|z, L) - \log q_\phi(z|x, L) + \log g(x, L) + \log p(L|z)], \qquad (5)$$

## 3.2 Training and Test

**Deterministic Context**. The lower-bound $\mathcal{Q}(x, L)$ transforms the intractable log-sum in Eq. (2) into the efficient sum-log in Eq. (3), which can be optimized by using Monte Carlo unbiased gradient estimator such as REINFORCE [50]. However, due to that $\phi$ is dependent on the sampling of $z$ over $\mathcal{O}(2^{|\mathcal{X}|})$ configurations, its gradient variance is large. To this end, we implement $q_\phi(z|x, L)$ as a differentiable but biased encoder:

$$z = f(x, L) = \sum_{x' \in \mathcal{X}} x' \cdot q_\phi(x'|x, L), \quad (6)$$

where we slightly abuse $q_\phi$ as a score function such that $\sum_{x'} q_\phi(x'|x, L) = 1$. Note that this deterministic context can be viewed as applying the "re-parameterization" trick as in Variational Auto-Encoder [19]: rewriting $z \sim q_\phi(z|x, L)$ to $z = f(x, L; \epsilon), \epsilon \sim p(\epsilon)$, where the stochasticity of the auxiliary random variable $\epsilon$ comes from training samples $x \in \mathcal{X}(\epsilon)$. A clear example is Adversarial Autoencoder [29] which shows that such stochasticity achieves similar test-likelihood compared to other distributions.

**Objective Function**. Applying Eq. (6) to Eq. (3), we can rewrite $Q(x, L)$ into a function of only one sample estimation, which is a common practice in SGD:

$$\mathcal{Q}(x, L) = \log p_\theta(x|z, L) - \log q_\phi(z|x, L) + \log p_\omega(z|L). \quad (7)$$

In supervised setting where the ground truth of the referent is known, to distinguish the referent from other objects, we need to train a model that outputs a high $p(x|L)$ (i.e., $\mathcal{Q}(x, L)$), while maintaining a low $p(x'|L)$ (i.e., $\mathcal{Q}(x', L)$), whenever $x' \neq x$. Therefore, we use the so-called Maximum Mutual Information loss as in [30] $-\log\{\mathcal{Q}(x, L)/\sum_{x'} \mathcal{Q}(x', L)\}$, where we do not need to explicitly model the distributions with normalizations; we use the following score function:

$$\mathcal{Q}(x, L) \propto \mathcal{S}(x, L) = s_\theta(x, L) - s_\phi(x, L) + s_\omega(x, L), \quad (8)$$

where $z$ is omitted as it is a function of $x$ in Eq. (6). $s_\theta$, $s_\phi$, and $s_\omega$ are the score functions (e.g., $p_\theta \propto s_\theta$) for $p_\theta$, $q_\phi$, and $p_\omega$, respectively. These functions will be detailed in Section 4.3. Similar to Eq. (8), we use the following score function to incorporate referring expression generation in the variational framework:

$$\mathcal{Q}'(x, L) \propto \mathcal{S}'(x, L) = s_\theta(x, L) - s_\phi(x, L) + s_{\omega'}(x, L) + s_\psi(x, L), \quad (9)$$

where $s_{\omega'}$ is the score function for $g(x, L)$, and shares the same structure with $s_\omega$; $s_\psi$ is the score function for $p(L|z)$.

In this way, maximizing Eq. (7) is equivalent to minimizing the following softmax loss:

$$\mathcal{L}_s = -\log \text{softmax}\, \mathcal{S}(x_{gt}, L), \quad (10)$$

where the softmax is over $x \in \mathcal{X}$ and $x_{gt}$ is the ground truth referent region.

Note that the reciprocity between referent and context can be extended to unsupervised learning, where neither of the referent and context has annotation. In this setting, we adopt the *image-level* max-pooled MIL loss functions for unsupervised referring expression grounding:

$$\mathcal{L}_u = -\max_{x \in \mathcal{X}} \log \text{softmax}\, \mathcal{S}(x, L), \quad (11)$$

where the softmax is over $x \in \mathcal{X}$. Note that the max-pooled MIL function is reasonable since there is only one ground truth referent given an expression and image training pair.

At test stage, in both supervised and unsupervised settings, we predict the referent region $x^*$ by selecting the region $x \in \mathcal{X}$ with the highest score:

$$x^* = \arg\max_{x \in \mathcal{X}} \mathcal{S}(x, L). \quad (12)$$

**Referring Expression Generation**. During training stage, instead of using the ground-truth referent, we sample a concrete referent region $\hat{x}$ from $p(x|L)$ and calculate its estimated context $\hat{z} = f(\hat{x}, L)$ using Eq. (6). The reason why we use the sampled referent is to punish false grounding results using the generation module. However, the generation loss $\mathcal{L}_G = \mathbb{E}_{x \sim p(x|L)} \mathcal{L}_c(x, L)$ is non-differentiable to referent grounding part using the concrete referent, where $\mathcal{L}_c(x, L)$ is the sum of cross entropy over the predicted words at each step. Hence, we apply policy gradient method in REINFORCE [50] for end-to-end training. The gradient $\nabla \mathcal{L}_G$ of the generation loss $\mathcal{L}_G$ is:

$$\nabla \mathcal{L}_G = E_{x \sim p(x|L)}[\mathcal{L}_c(x, L)\nabla \log p(x|L) + \nabla \mathcal{L}_c(x, L)]. \quad (13)$$

In practice, the gradient $\nabla \mathcal{L}_G$ can be estimated using Monte-Carlo sampling as:

$$\nabla \mathcal{L}_G \approx \frac{1}{K} \sum_{k=1}^N [\mathcal{L}_c(x_k, L)\nabla \log p(x_k|L) + \nabla \mathcal{L}_c(x_k, L)], \quad (14)$$

where $x_k$ is sampled from $p(x|L)$. We simply use $K = 1$ in our implementation. Since $p(x|L)$ is fully differentiable, the gradient can be transferred to referent grounding part via backpropagation. Following [48], we utilize a moving average baseline $b$ to reduce the variance of estimated gradient using REINFORCE, and replace $\mathcal{L}_c(x_k, L)$ with $\mathcal{L}_c(x_k, L) - b$ in Eq. (14). The baseline $b_t$ at $t$-th iteration is estimated by accumulating the previous losses $\mathcal{L}_c(x, L)$ with exponential decay:

$$b_t = 0.9 \times b_{t-1} + 0.1 \times \mathcal{L}_c(x_{k_t}, L). \quad (15)$$

## 4 MODEL ARCHITECTURE

The overall architecture of the proposed variational context framework is illustrated in Fig. 2. Thanks to the deterministic context in Eq. (6) and REINFORCE in Eq. (14), the six modules in our model can be integrated into an end-to-end differentiable fashion. Next, we will detail the implementation of each module.

### 4.1 RoI Features

Given an image with a set of Region of Interests (RoIs) $\mathcal{X}$, obtained by any off-the-shelf proposal generator [64] or object detectors [25], this module extracts the feature vector $\mathbf{x}_i$ for every RoI. In particular, $\mathbf{x}_i$ is the concatenation of visual feature $\mathbf{v}_i$ and spatial feature $\mathbf{p}_i$. For $\mathbf{v}_i$, we can use the output of a pre-trained convolutional network (cf. Section 5). If the object category of each RoI is available, we can further utilize the comparison between the referent and other objects to capture the visual difference such as "the largest/baby elephant". Specifically, we append the visual difference (visdif) feature [57] $\delta\mathbf{v}_i = \frac{1}{n}\sum_{j \neq i} \frac{\mathbf{v}_i - \mathbf{v}_j}{||\mathbf{v}_i - \mathbf{v}_j||}$ to the
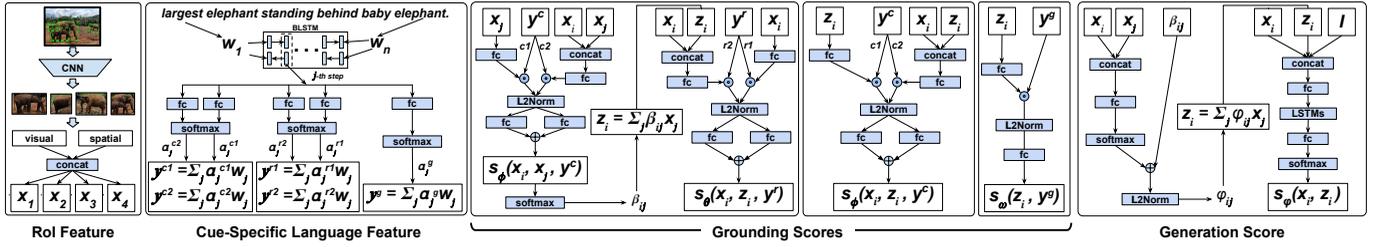
Fig. 2. The architecture of the proposed Variational Context framework. It consists of a region feature extraction module (Section 4.1), and a language feature extraction module (Section 4.2), three grounding modules (Section 4.3), and one generation module (Section 4.4). It can be trained in an end-to-end fashion with the input of a set of image regions and a referring expression, using the supervised loss (Eq. (10)) or the unsupervised loss (Eq. (11)). LSTMs: Long short-term memory networks. fc: fully-connected layer. concat: vector concatenation. L2Norm: L2 normalization layer. $\odot$: element-wise vector multiplication. $\oplus$: add.
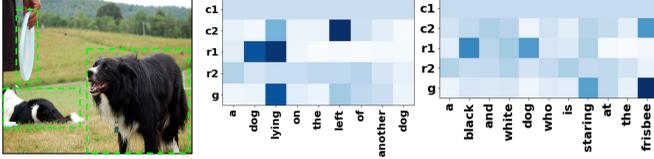


Fig. 3. Two qualitative examples of the cue-specific language feature word weights. Darker color indicates higher weights. c/r+1/2: context/referent-cue + single/pairwise.

original $\mathbf{v}_i$ visual feature, where $n$ is the number of objects chosen for comparison (*e.g.*, the number of RoI in the same object category). For spatial feature, we use the 5-d spatial attributes $\mathbf{p}_i = \left[\frac{x_{tl}}{W}, \frac{y_{tl}}{H}, \frac{x_{br}}{W}, \frac{y_{br}}{H}, \frac{w \cdot h}{W \cdot H}\right]$, where $x$ and $y$ are the coordinates the top left (tl) and bottom right (br) RoI of the size $w \times h$, and the image is of the size $W \times H$.

## 4.2 Cue-Specific Language Features

An alternative for language understanding is to employ external NLP parsers to parse the referring expression into several textual components, such as subject, object, and relationship. However, conventional parsers (*e.g.*, Standford Dependency) are observed to be suboptimal to the grounding referring expression task [15], and cannot be end-to-end trained. In this work, we choose to parse the referring expression with language attention by learning in an end-to-end fashion. Specifically, the referring expression is represented into different cue-specific language features, which is inspired by the attention weighted sum of word vectors [3], [15], [27]. We parameter the weights by context-cue $\phi$, referent-cue $\theta$, and generic-cue $\omega$ based on their different purposes: the context-cue feature helps to estimate context prior, the referent-cue feature aims to localize the referent, and the generic-cue feature encourages the estimated context to be consistent with the visual context described in the expression. Formally, the context-cue language feature $\mathbf{y}^c = [\mathbf{y}^{c1}, \mathbf{y}^{c2}]$ is a concatenation of $\mathbf{y}^{c1}$ for language-vision association between *single* RoI and the expression, and $\mathbf{y}^{c2}$ for the association between *pairwise* RoIs; the referent-cue language feature $\mathbf{y}^r$ can be represented in a similar way to $\mathbf{y}^c$; the generic-cue language feature $\mathbf{y}^g$ is only for single RoI association as it is an independent prior. The weights of each cue are calculated from the hidden state vectors of a 2-layer bi-directional LSTM (BLSTM) [40], scanning through the expression. The

hidden states encode forward and backward compositional semantic meanings of the sentences, beneficial for selecting words that are useful for single and pairwise associations. Specifically, suppose $\mathbf{h}_j$ as the 4,000-d concatenation of forward and backward hidden vectors of the $j$-th word, without loss of generality, the word attention weight $\alpha_j$ and the language feature $\mathbf{y}$ for single/pairwise association of any cue can be calculated as:

$$\mathbf{m}_j = \mathrm{fc}(\mathbf{h}_j), \alpha_j = \mathrm{softmax}_j(\mathbf{m}_j), \mathbf{y} = \sum_j \alpha_j \mathbf{w}_j, \quad (16)$$

where $\mathbf{w}_j$ is a 300-d vector. Note that the BLSTM module can be jointly trained with the entire model.

Fig. 3 shows that the cue-specific language features dynamically weight words in different expressions. We can have two interesting observations. First, c1 is almost uniform while c2 is highly skewed; although r2 is more skewed than c1, it is still less skewed than r1. This is reasonable since: 1) without ground-truth, individual score (c1) does not help much for context estimation from scratch; context is more easily found by the pairwise score (c2) induced by relationships or other objects (*e.g.*, "left" or "frisbee"); 2) in referent grounding with ground truth, individual score (r1) is sufficient (*e.g.*, "dog lying" and "black white dog") and pairwise score (r2) is helpful; 3) g is adaptive to the number of object categories in the expression, *i.e.*, if the context object is of the same category as the referent, g weighs descriptive or relationship words higher (*e.g.*, "lying, standing, left"), and nouns higher (*e.g.*, "frisbee"), otherwise; moreover, it demonstrates that the deterministic guess of $z$ in Eq. (6) is meaningful.

## 4.3 Score Functions for Comprehension

For any image and expression pair, given the RoI feature $\mathbf{x}_i$, and the cue-specific language feature $\mathbf{y}^c$, $\mathbf{y}^r$, and $\mathbf{y}^g$, we implement the final grounding score in Eq. (8) as:

$$\begin{aligned} \mathbf{z}_i = \sum_j \mathrm{softmax}_j \left( s_\phi(\mathbf{x}_i, \mathbf{x}_j, \mathbf{y}^c) \right) \mathbf{x}_j, \\ s_\theta(x, L) \leftarrow s_\theta(\mathbf{x}_i, \mathbf{z}_i, \mathbf{y}^r), \\ s_\phi(x, L) \leftarrow s_\phi(\mathbf{x}_i, \mathbf{z}_i, \mathbf{y}^c), \\ s_\omega(x, L) \leftarrow s_\omega(\mathbf{z}_i, \mathbf{y}^g), \end{aligned} \quad (17)$$

where the right-hand side functions are defined as below.

**Context Estimation Score**: $s_\phi(\mathbf{x}_i, \mathbf{x}_j, \mathbf{y}^c)$. It is a score function for modeling the context posterior $q_\phi(z|x, L)$, *i.e.*, given an RoI $\mathbf{x}_i$ as the candidate referent, we calculate the

likelihood of any RoI $\mathbf{x}_j$ to be the context. We can also use this function to estimate the final context posterior score $s_\phi(\mathbf{x}_i, \mathbf{z}_i, \mathbf{y}^c)$. Specifically, the context estimation score is a sum of the single and pairwise vision-language association scores: $\mathbf{x}_j$ and $\mathbf{y}^{c1}$, $[\mathbf{x}_i, \mathbf{x}_j]$ and $\mathbf{y}^{c2}$. Each associate score is an fc output from the input of a normalized feature:

$$
\begin{aligned}
\mathbf{m}_j^1 &= \mathbf{y}^{c1} \odot \mathrm{fc}(\mathbf{x}_j), \quad \mathbf{m}_j^2 = \mathbf{y}^{c2} \odot \mathrm{fc}([\mathbf{x}_i, \mathbf{x}_j]), \\
\widetilde{\mathbf{m}}_j^1 &= \mathrm{L2Norm}(\mathbf{m}_j^1), \quad \widetilde{\mathbf{m}}_j^2 = \mathrm{L2Norm}(\mathbf{m}_j^2), \qquad (18) \\
s_\phi(\mathbf{x}_i, \mathbf{x}_j, \mathbf{y}^c) &= \mathrm{fc}(\widetilde{\mathbf{m}}_j^1) + \mathrm{fc}(\widetilde{\mathbf{m}}_j^2),
\end{aligned}
$$

where the element-wise multiplication $\odot$ is an effective way for fusing multimodal features [2]. According to Eq. (6), we can obtain the estimated context $z$ as $\mathbf{z}_i = \sum_j \beta_j \mathbf{x}_j$, where $\beta_j = \mathrm{softmax}_j(s_\phi(\mathbf{x}_i, \mathbf{x}_j, \mathbf{y}^c))$.

**Referent Grounding Score**: $s_\theta(\mathbf{x}_i, \mathbf{z}_i, \mathbf{y}^r)$. After obtaining the context feature $\mathbf{z}_i$, we use this score function to calculate how likely a candidate RoI $\mathbf{x}_i$ is the referent given the context $\mathbf{z}_i$. This function is similar to Eq. (18).

**Context Regularization Score**: $s_\omega(\mathbf{z}_i, \mathbf{y}^g) - s_\phi(\mathbf{x}_i, \mathbf{z}_i, \mathbf{y}^c)$. As discussed in Eq. (8), this function scores how likely the estimated context feature $\mathbf{z}_i$ is consistent with the content mentioned in the expression. In particular, $s_\omega(\mathbf{z}_i, \mathbf{y}^g)$ is only dependent on single RoI:

$$
\mathbf{m}_i = \mathbf{y}_i^g \odot \mathrm{fc}(\mathbf{z}_i), \widetilde{\mathbf{m}}_i = \mathrm{L2Norm}(\mathbf{m}_i), s_\omega(\mathbf{z}_i, \mathbf{y}_i^g) = \mathrm{fc}(\mathbf{m}_i). \quad (19)
$$

### 4.4 Score Function for Generation

For any referent and expression pair, given the RoI feature $\mathbf{x}_i$ and the context-specific language feature $\mathbf{y}^c$, we implement the generation score function in Eq. (9) to reconstruct the referring expression as:

$$
\begin{aligned}
\hat{\mathbf{z}}_i &= \sum_j \varphi_j \mathbf{x}_j, \\
s_\psi(x, L) &\leftarrow s_\psi(\mathbf{x}_i, \hat{\mathbf{z}}_i),
\end{aligned} \qquad (20)
$$

where $\hat{\mathbf{z}}_i$ represents the estimated context for generation, and the joint region attention $\varphi_j$ is defined as below:

$$
\begin{aligned}
\beta_j &= \mathrm{softmax}_j(s_\phi(\mathbf{x}_i, \mathbf{x}_j, \mathbf{y}^c)), \\
\gamma_j &= \mathrm{softmax}_j(\mathrm{fc}([\mathbf{x}_i, \mathbf{x}_j])), \qquad (21) \\
\varphi_j &= \mathrm{L2Norm}_j(\beta_j \odot \gamma_j),
\end{aligned}
$$

where $\beta_j$ and $\gamma_j$ represents the region attention weight for comprehension and generation respectively, and $\varphi_j$ mixes these attention weights using the element-wise multiplication $\odot$. Note that $\beta_j$ shares the same calculation with context estimation score in the comprehension module, which can evaluate the estimated context. We then modify the vanilla language generation model [47] to reconstruct the referring expression. Different from [47], we feed the referent $\mathbf{x}_i$ with its context $\hat{\mathbf{z}}_i$ and the image feature $\mathbf{I}$ into the LSTM model for sequence generation:

$$
\begin{aligned}
\mathbf{w}_{-1} &= \mathrm{fc}([\mathbf{x}_i, \hat{\mathbf{z}}_i, \mathbf{I}]), \quad \mathbf{h}_{-2} = \mathbf{0}, \\
\mathbf{w}_t &= W_e S_t, \quad \mathbf{h}_t = \mathrm{LSTM}(\mathbf{w}_t, \mathbf{h}_{t-1}), \\
\mathbf{p}_t &= \mathrm{softmax}(\mathrm{fc}(\mathbf{h}_t)), \qquad (22) \\
s_\psi(\mathbf{x}_i, \hat{\mathbf{z}}_i) &= \prod_t \mathbf{p}_t^T S_{t+1}
\end{aligned}
$$

where $W_e$ is the word embedding matrix shared with the comprehension module, and $S_t$ is the one-hot encoding for the input word $w_t$ at step $t$. Note that the start word and stop word are denoted by $w_0$ and $w_{T+1}$, respectively, standing for the beginning and end of the referring expression, where $T$ is the length of the referring expression. A complete referring expression is generated when the LSTM encounters the stop word or the length of expression reaches the maximum number.

## 5 EXPERIMENT

### 5.1 Datasets

We used four popular benchmarks for the referring expression grounding task.

**RefCOCO** [57]. It has 142,210 referring expressions for 50,000 referents (*e.g.*, object instances) in 19,994 images from MSCOCO [23]. The expressions are collected in an interactive way [18]. The dataset is split into train, validation, Test A, and Test B, which has 120,624, 10,834, 5,657 and 5,095 expression-referent pairs, respectively. An image contains multiple people in Test A and multiple objects in Test B.

**RefCOCO+** [57]. It has 141,564 expressions for 49,856 referents in 19,992 images from MSCOCO. The difference from RefCOCO is that it only allows appearances but no locations to describe the referents. The split is 120,191, 10,758, 5,726 and 4,889 expression-referent pairs for train, validation, Test A, and Test B respectively.

**RefCOCOg** [30], [56]. It has 95,010 referring expressions for 49,822 objects in 25,799 images from MSCOCO. Different from RefCOCO and RefCOCO+, this dataset not collected in an interactive way and contains longer sentences containing both appearance and location expressions. RefCOCOg has two types of split. The old split [30] is 85,474 and 9,536 expression-referent pairs for training and validation. It should be noticed that the old partitioned by objects, thus some images may exist in both train and validation sets. We represent the validation set as "Val*". The new partition [56] randomly divides images into train, validation and test set. The split is 80,512, 4,896 and 9,602 expression-referent pairs for train, validation and test, respectively. The validation and test sets are represented as "Val" and "Test". Compared to RefCOCO and RefCOCO+, RefCOCOg contains longer expressions, which makes it more challenging for comprehension and generation.

**RefCLEF** [18]. It contains 20,000 images with annotated image regions. It has some ambiguous (e.g. anywhere) phrases and mistakenly annotated image regions that are not described in the expressions. For fair comparison, we used the split released by [16], [39], *i.e.*, 58,838, 6,333 and 65,193 expression-referent pairs for training, validation and test, respectively.

### 5.2 Settings and Metrics

For comprehension module, we used an English vocabulary of 72,704 words contained in the GloVe pre-trained word vectors [34], which was also used for the initialization of our word vectors. We used a "unk" symbol for the input word of the BLSTM if the word is out of the vocabulary; we set the sentence length to 20 and used "pad" symbol to pad expression sentence $< 20$. For RoI visual features on RefCOCO, RefCOCO+, and RefCOCOg which
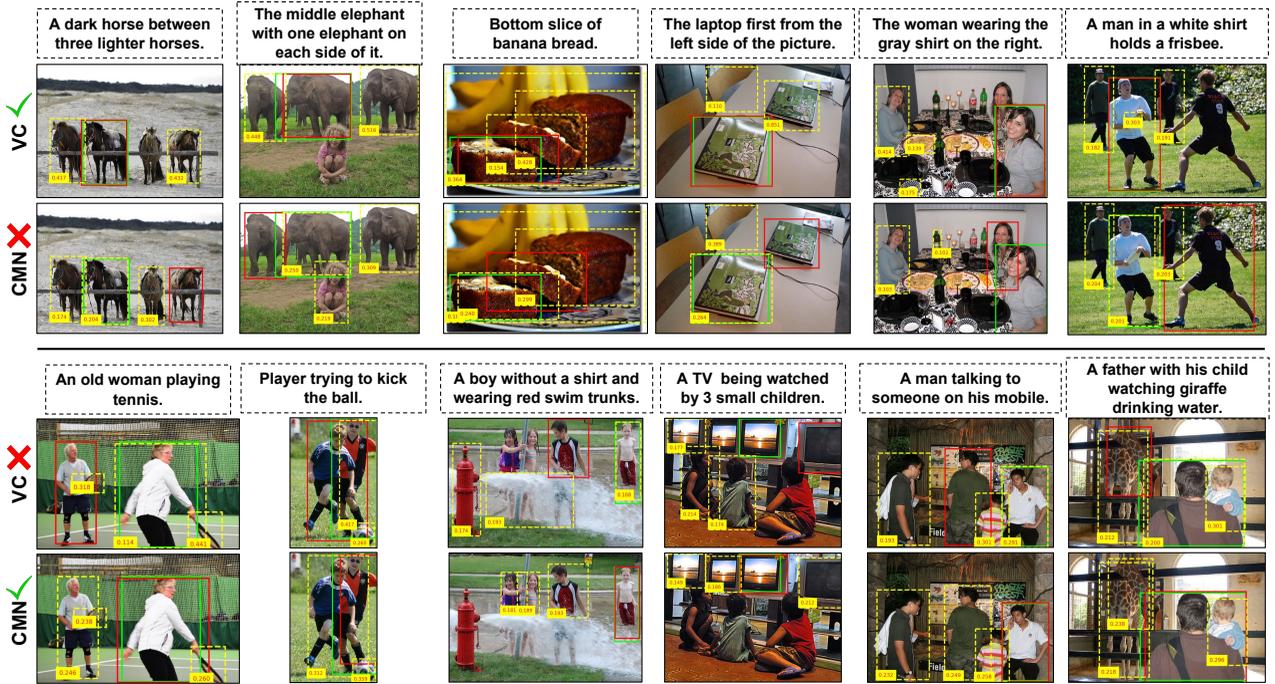
Fig. 4. Qualitative results on RefCOCOg (det) showing comparisons between correct (green tick) and wrong referent grounds (red cross) by VC and CMN using VGG features. The denotations of the bounding box colors are as follows. Solid red: grounding referent; solid green: ground truth; dashed yellow: grounding context. We only display top 3 context objects with the context ground probability $> 0.1$. We can observe that VC has more reasonable context localizations than CMN, even in cases when the referent ground of VC fails.


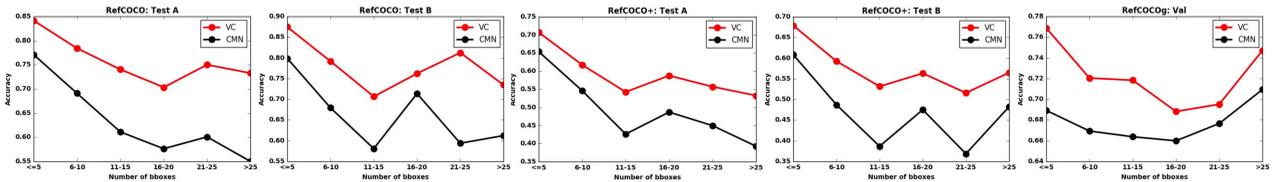
Fig. 5. Performances of VC and CMN with different number of object bounding boxes on RefCOCO Test A &B, RefCOCO+ Test A & B, and RefCOCOg Val. Compared to CMN, we can see that VC is more effective in context modeling when the number of objects is large.

have MSCOCO annotated regions with object categories, we used the concatenation of the 4,096-d fc7 output of a VGG-16 based Faster-RCNN network [38] trained on MSCOCO and its corresponding 4,096-d visdif feature [57]; although RefCLEF regions also have object categories, for fair comparison with [39], we did not use the visdif feature. For generation module, we built an additional vocabulary including words which occur at least 5 times in the training set. The maximum length of generated sentences is set as 20. The hidden state size of LSTM is set as 512. We also regularized the LSTM using dropout with ratio of 0.3. Following [14], we also use an entropy regularization $5 \times 10^{-3}$ over the conditional distribution $p(x|L)$ to encourage exploration through the sampling space. For fair comparison with [56], we also used the average-pooled C4 feature and phrase-guided embedding, which are provided by [56].

The model training is single-image based, with all referring expressions annotated. We applied SGD of 0.95-momentum with initial learning rate of 0.01, multiplied by 0.1 after every 120,000 iterations, up to 160,000 iterations. Parameters in BILSTM and fc-layers were initialized by Xavier [12] with $5 \times 10^{-4}$ weight decay. Other settings

were default in TensorFlow. Note that our model is trained without bells and whistles, therefore, other optimization tricks such as batch normalization [17] and GRU [5] are expected to further improve the results reported here. Besides the ground truth annotations, grounding to automatically detected objects is a more practical setting. Therefore, we also evaluated with detected objects, the SSD-detected bounding boxes [25] provided by [58] using VGG-based model, and Faster R-CNN detected bounding boxes provided by [56] using ResNet-based model. A grounding is considered as correct if the intersection-over-union (IoU) of the top-1 scored region and the ground-truth object is larger than 0.5. The grounding accuracy (a.k.a, P@1) is the fraction of correctly grounded test expressions.

### 5.3 Evaluations of Supervised Grounding

We compared our variational context (VC) method with state-of-the-art referring expression methods published in recent years, which can be categorized into: 1) generation-comprehension based such as MMI [30], Attr [24], Speaker [58], Listener [58], and SCRC [16]; 2) localization based such as GroundR [39], NegBag [33],
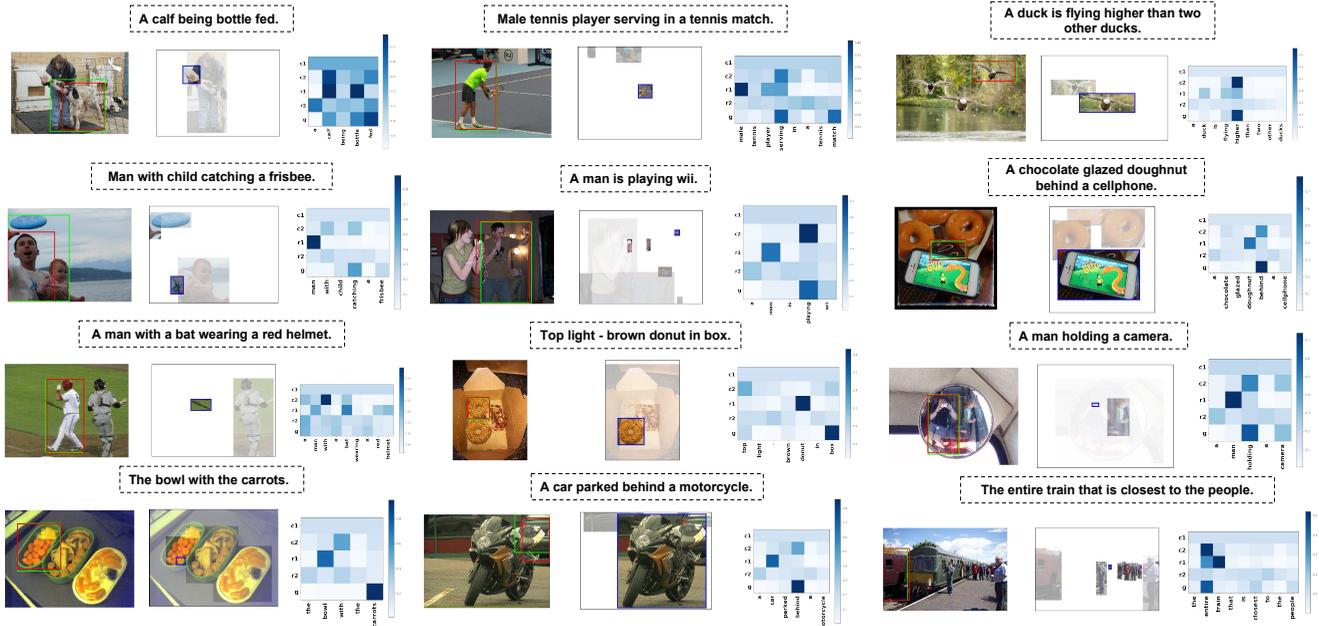
Fig. 6. Qualitative results of our full model (VC w/ Gen+PG) on RefCOCOg (det). The first column shows the grounding results. The second column shows the context estimation results. The third column shows the cue-specific language feature word weights. The denotations of the bounding box colors are as follows. Solid red: grounding referent; solid green: ground truth; solid blue: grounding context with highest probability.

TABLE 1
Supervised grounding performances (Acc%) of comparing methods using VGG features on MSCOCO ground-truth regions. Note that [58] reports slightly higher accuracies using ensemble models of Listener and Speaker. For fair comparison, we only report their single models.

| Dataset | RefCOCO | | RefCOCO+ | | RefCOCOg |
|---|---|---|---|---|---|
| Split | Test A | Test B | Test A | Test B | Val* |
| MMI [30] | 71.72 | 71.09 | 58.42 | 51.23 | 62.14 |
| NegBag [33] | 75.6 | 78.0 | — | — | 68.4 |
| Attr [24] | 78.85 | 78.07 | 61.47 | 57.22 | 69.83 |
| CMN [15] | 75.94 | 79.57 | 59.29 | 59.34 | 69.30 |
| Speaker [58] | 78.95 | 80.22 | 64.60 | 59.62 | 72.63 |
| Listener [58] | 78.45 | 80.10 | 63.34 | 58.91 | 72.25 |
| PLAN [63] | 80.81 | 81.32 | 66.31 | 61.46 | 69.47 |
| A-ATT [8] | **81.17** | 80.01 | **68.76** | 60.63 | 73.18 |
| MAttNet [56] | 79.99 | 82.30 | 65.04 | 61.77 | 73.08 |
| VC w/o reg | 75.59 | 79.69 | 60.76 | 60.14 | 71.05 |
| VC w/o $\alpha$ | 74.03 | 78.27 | 57.61 | 54.37 | 65.13 |
| VC | 78.98 | **82.39** | 62.56 | **62.90** | **73.98** |

TABLE 2
Supervised grounding performances (Acc%) of comparing methods using VGG features on MSCOCO detected regions. Note that [58] reports slightly higher accuracies using ensemble models of Listener and Speaker. For fair comparison, we only report their single models.

| Dataset | RefCOCO | | RefCOCO+ | | RefCOCOg |
|---|---|---|---|---|---|
| Split | Test A | Test B | Test A | Test B | Val* |
| MMI [30] | 64.90 | 54.51 | 54.03 | 42.81 | 45.85 |
| NegBag [33] | 58.6 | 56.4 | — | — | 39.5 |
| Attr [24] | 72.08 | 57.29 | 57.97 | 46.20 | 52.35 |
| CMN [15] | 71.03 | 65.77 | 54.32 | 47.76 | 57.47 |
| Speaker [58] | 72.95 | 63.43 | **60.43** | 48.74 | 59.51 |
| Listener [58] | 72.95 | 62.98 | 59.61 | 48.44 | 58.32 |
| PLAN [63] | **75.31** | 65.52 | 61.34 | 50.86 | 58.03 |
| VC w/o reg | 70.78 | 65.10 | 56.82 | 51.30 | 60.95 |
| VC w/o $\alpha$ | 70.73 | 64.63 | 53.33 | 46.88 | 55.72 |
| VC | 73.33 | **67.44** | 58.40 | **53.18** | **62.30** |

CMN [15], MAttNet [56], PLAN [63], A-ATT [8]. Note that NegBag and CMN are MIL-based (multiple instance learning) models. In particular, we used the author-released code to obtain the results of CMN on RefCLEF, RefCOCO, and RefCOCO+.

**Single comprehension module.** From the results of VGG-based models on RefCOCO, RefCOCO+, and RefCOCOg in Table 1 and 2, and that on RefCLEF in Table 3, we can see that VC achieves the state-of-the-art performance. We believe that the improvement is attributed to the variational Bayesian modeling of context. First, on all datasets, except for the most recent reinforcement learning based method [58] or multiple attention mechanism based mathods [8], [56], [63] on the Test A split, VC outperforms all the other sentence generation-comprehension methods that do not model context. Second, compared to VC without

the regularization term in Eq. (3) (VC w/o reg), VC can boost the performance by around 2% on all datasets. This demonstrates the effectiveness of the KL divergence for the prevention of the overfitted context estimation.

In particular, we further demonstrate the superiority of VC over the most recent MIL-based method CMN. As illustrated in Fig. 4, VC has better context comprehension in both of the language and image regions than CMN. For example, in the top two rows where VC is correct and CMN is wrong, for the grounding in the second column, CMN unnecessarily considers the "girl" as context but the expression only describes using "elephant"; in the last column, CMN misses the key context "frisbee". Even in the failure cases where VC is wrong and CMN is correct, VC still localizes reasonable context. For example, in the fourth column, although CMN grounds the correct TV, it is based on incorrect context of other TVs; while VC can predict the
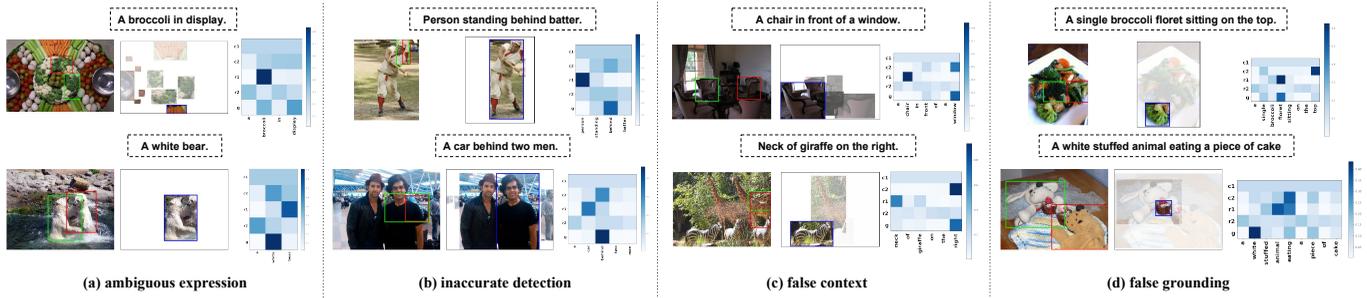
Fig. 7. Common failure cases of our full model in supervised grounding on RefCOCOg. Each example shows grounding results, context estimation results, and cue-specific features from left to right. The denotations of the bounding box colors are as follows. Solid red: grounding referent; solid green: ground truth; solid blue: grounding context with highest probability.

TABLE 3
Performances (Acc%) of supervised and unsupervised methods on RefCLEF.

|  | Sup. | Sup. (det) | Unsup. (det) |
|---|---|---|---|
| SCRC [16] | 72.74 | 17.93 | — |
| GroundR [39] | — | 26.93 | 10.70 |
| CMN [15] | 81.52 | 28.33 | — |
| VC | **82.43** | **31.13** | 14.11 |
| VC w/o $\alpha$ | 79.60 | 27.40 | **14.50** |

TABLE 4
Supervised grounding performances (Acc%) of ablation study using generation module or better visual representation on MSCOCO ground-truth regions. † and ‡ indicates that this model uses res101 feature and attribute-based phrase-guided feature [56], respectively.

|  | RefCOCO | | RefCOCO+ | | RefCOCOg | | |
|---|---|---|---|---|---|---|---|
|  | Test A | Test B | Test A | Test B | Val* | Val | Test |
| MAttNet [56]† | 81.58 | 83.34 | 66.59 | 65.08 | — | 75.96 | 74.56 |
| MAttNet [56]†‡ | 85.26 | 84.57 | 75.13 | 66.17 | — | 78.10 | 78.12 |
| VC | 78.98 | 82.39 | 62.56 | 62.90 | 73.98 | 74.61 | 74.58 |
| VC w/ Gen | 79.16 | 82.04 | 62.84 | 62.88 | 74.20 | 74.98 | 75.06 |
| VC w/ Gen+PG | 79.30 | 82.04 | 63.22 | 63.12 | **74.96** | 75.35 | 75.11 |
| VC w/ Gen+PG† | 80.40 | 83.51 | 67.52 | 66.46 | — | 77.49 | 76.64 |
| VC w/ Gen+PG†‡ | **86.26** | **85.00** | **76.48** | **68.13** | — | **79.80** | **79.96** |

correct context "children". In addition, we observed that most of the cases that CMN is better than VC involves multiple humans. This demonstrates that VC is better at grounding objects of different categories. VC is also effective in images with more objects. Fig. 5 shows the performances of VC and CMN with various number of bounding boxes. We can observe that VC considerably outperforms CMN over all bounding boxes numbers. Recall that context is the key to distinguish objects of the same category. In particular, on the Test A sets of RefCOCO and RefCOCO+ where the grounding is only about people, *i.e.*, the same object category, the gap between VC and CMN is becoming larger as the box number increases. This demonstrates that MIL is ineffective in modeling context, especially when the number of image regions is large.

**Cooperation with generation module.** Furthermore, we exploit the grounding performance of VC incorporating referring expression generation in the variational framework using Eq. (9). As shown in Table 4, VC incorporating generation module using the policy gradient method REINFORCE (VC w/ Gen+PG) achieves the best performance except the Test B split of RefCOCO, which means that referring expression generation can help with comprehension in our framework. Note that we only use the single comprehension module during test time, which has the same structure with better-learned parameters compared to the single VC model. Note that VC w/ Gen slightly performs worse than VC on two splits. The possible reason comes from the difficulty of language understanding. The average length of referring expression in RefCOCO and RefCOCO+ is about 3.6, while the queries in RefCOCOg have an average length of around 8.4. This observation indicates that context estimation plays a more important role in grounding long descriptions than short phrases, since long description tends to include more complex context information.

**Better visual representation.** Recently, [56] use ResNet-FPN backbone for feature extraction instead of VGG. We also evaluate VC model using ResNet-based Faster-RCNN visual representation, which can further improve the grounding performance by at least 1%, especially on RefCOCO+ dataset. In additional to ResNet feature, MAttNet also includes attention mechanism to obtain phrase-guided embedding in cooperation with object attribute prediction. For fair comparison, we just use their pre-trained model to extract attribute-based phrase-guided feature, and concatenate it with origin visual feature. As shown in Table 4, VC w/ Gen+PG yields the state-of-the-art MAttNet by at least 1% on RefCOCO+ and RefCOCOg datasets. It is worth noting that the relationship module in MAttNet assumes that only one object contributes to the context, which suffers from the ineffectiveness of MIL in modeling context.

**Qualitative results and failure cases studies.** The qualitative results of our full model on RefCOCOg dataset are shown in Fig. 6. As illustrated in Fig. 6, our full model estimates reasonable context and cue-specific feature. For example, for the referring expression "a calf being bottle fed", the word "calf" is the key to referent-cue feature. Since there are two calves in the image, our full model highlights the word "fed" in context-cue and generic-cue features. Meanwhile, our full model correctly focuses on the "fed" relation represented by the visual content of bottle and man in context estimation. Some common failure cases on RefCOCOg are illustrated in Figure 7, classified into four groups: ambiguous expression, inaccurate detection, false context, and false grounding. Ambiguous expression means that the expression matches more than one candidate objects. Actually the grounding results in first row are correct since they also correspond to the referring expression. Inaccurate detection means that
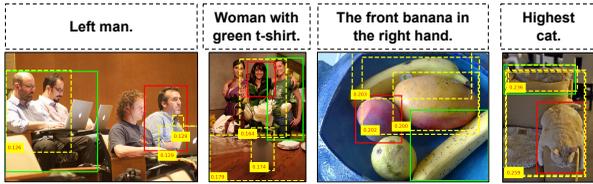
Fig. 8. Common failure cases in unsupervised grounding with detected bounding boxes. From left to right: RefCOCO, RefCOCO+, and RefCOCOg. The failure is mainly to the challenging unsupervised relation modeling between referent and context.

TABLE 5
Unsupervised grounding performances (Acc%) of comparing methods using VGG features on RefCOCO, RefCOCO+, and RefCOCOg.

| Dataset | RefCOCO | | RefCOCO+ | | RefCOCOg |
|---|---|---|---|---|---|
| Split | Test A | Test B | Test A | Test B | Val* |
| VC w/o reg | 13.59 | 21.65 | 18.79 | 24.14 | 25.14 |
| VC | 17.34 | 20.98 | 23.24 | 24.91 | **33.79** |
| VC w/o $\alpha$ | **33.29** | **30.13** | **34.60** | **31.58** | 30.26 |
| Dataset | RefCOCO (det) | | RefCOCO+ (det) | | RefCOCOg (det) |
| Split | Test A | Test B | Test A | Test B | Val* |
| VC w/o reg | 17.14 | 22.30 | 19.74 | 24.05 | 28.14 |
| VC | 20.91 | 21.77 | 25.79 | 25.54 | **33.66** |
| VC w/o $\alpha$ | **32.68** | **27.22** | **34.68** | **28.10** | 29.65 |

the referent is partially or even not detected due to the limitation of detector. Although our model successfully localizes the referent, only a part of target object has been detected. False context means that the context estimation result is incorrect. The failure estimated context confuses the grounding module to distinguish the ground-truth referent from other objects of the same category. False grounding means that the grounding module fails although the detection and estimated context are correct.

## 5.4 Evaluations of Unsupervised Grounding

We follow the unsupervised setting in GroundR [39]. To our best knowledge, it is the only work on unsupervised referring expression grounding. Note that it is also known as "weakly supervised" detection [60] as there is still image-level ground truth (*i.e.*, the referring expression). Table 3 reports the unsupervised results on the RefCLEF. We can see that VC outperforms the state-of-the-art GroundR, which is a generation-comprehension based method. This demonstrates that using context also helps unsupervised grounding. As there is no published unsupervised results on RefCOCO, RefCOCO+, and RefCOCOg, we only compared our baselines on them in Table 5. We can have the following three key observations which highlight the challenges of unsupervised grounding:

**Context Prior**. VC w/o reg is the baseline without the KL divergence as a context regularization in Eq. (3). We can see that in most of the cases, VC considerably outperforms VC w/o reg by over 2%, even over 5% on RefCOCO+ (det) and RefCOCOg (det). Note that this improvement is significantly higher than that in supervised setting (*e.g.*, $< 3\%$ as reported in Table 1). The reason is that the context estimation in Eq. (6) would be easier to be stuck in image regions that are irrelevant to the expression in unsupervised setting, therefore, context prior is necessary.
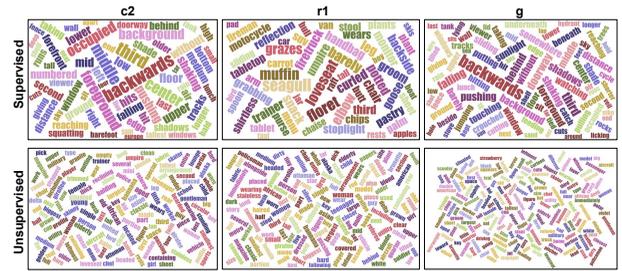


Fig. 9. Word cloud visualizations of cue-specific word attention $\alpha$ in Eq. (16) of context-cue (c2), referent-cue (r1), and generic-cue (g) using supervised (top row) and unsupervised training (bottom row) on RefCOCOg. Without supervision, it is difficult to discover meaningful language compositions.

**Language Feature**. Except on RefCOCOg, we consistently observed the *ineffectiveness* of the cue-specific language feature in unsupervised setting, *i.e.*, VC w/o $\alpha$ outperforms VC in Table 3 and 5. Here $\alpha$ represents the cue-specific word attention. This is contrary to the observation in the supervised setting as listed in Table 1 and 2, where VC w/o $\alpha$ is consistently lower than VC. Note that without the cue-specific word attention $\alpha$ in Eq. (16), the language feature is merely the average value of the word embedding vectors in the expression. In this way, VC w/o $\alpha$ does not encode any structural language composition as illustrated in Fig. 3, thus, it is better for short expressions. However, when the expression is long in RefCOCOg, discarding the language structure still degrades the performance on RefCOCOg.

**Unsupervised Relation Discovery**. Although we demonstrated that VC improves the unsupervised grounding by modeling context, we believe that there is still a large space for improving the quality of modeling the context. As the failure examples shown in Fig. 8, 1) many context estimations are still out of the scope of the expression, *e.g.*, we may localize the "cup" and "table" as context even though the expression is "woman with green t-shirt"; 2) we may mistake due to the wrong comprehension of the relations, *e.g.*, "right" as "left", even if the objects belong to the same category, *e.g.*, "elephant". For further investigation, Fig. 9 visualizes the cue-specific word attentions in supervised and unsupervised settings. The almost identical word attentions in unsupervised setting reflect the fact that the relation modeling between referent and context is not as successful as in supervised setting. This inspires us to exploit stronger prior knowledge such as language structure [52] and spatial configurations [49], [60].

## 5.5 Evaluation of Generation

For the generation task, we first evaluate our models using BLEU, METEOR and CIDEr, which are widely used evaluation metrics in generated description evaluation. The automatic evaluation results using above metrics are given in Table 6. Here, "Gen" represents the generation module trained separately without comprehension loss. From Table 6, we observe that the comprehension module helps to improve all the metrics significantly compared to the single generation module. This observation indicates that the estimated comprehension context can help to promote the performance of the generation module. In

TABLE 6
Automatic metrics on referring expression generation. Note that [58] reports slightly higher accuracy using reranking mechanism. For fair comparison, we only report their performance without reranking.

| | RefCOCO (Test A) | | | |
|---|---|---|---|---|
| | BLEU-1 | BLEU-2 | METEOR | CIDEr |
| MMI [30] | 0.478 | 0.295 | 0.175 | - |
| visdif [57] | 0.505 | 0.322 | 0.184 | - |
| Speaker [58] | - | - | **0.268** | 0.697 |
| Gen | 0.472 | 0.299 | 0.170 | 0.641 |
| VC w/ Gen | 0.548 | 0.361 | 0.188 | 0.707 |
| VC w/ Gen+PG | **0.556** | **0.368** | 0.194 | **0.716** |
| | RefCOCO (Test B) | | | |
| | BLEU-1 | BLEU-2 | METEOR | CIDEr |
| MMI [30] | 0.547 | 0.341 | 0.228 | - |
| visdif [57] | 0.583 | 0.382 | 0.245 | - |
| Speaker [58] | - | - | **0.329** | 1.323 |
| Gen | 0.548 | 0.351 | 0.237 | 1.271 |
| VC w/ Gen | 0.628 | 0.424 | 0.245 | 1.356 |
| VC w/ Gen+PG | **0.639** | **0.430** | 0.252 | **1.364** |
| | RefCOCO+ (Test A) | | | |
| | BLEU-1 | BLEU-2 | METEOR | CIDEr |
| MMI [30] | 0.370 | 0.203 | 0.136 | - |
| visdif [57] | 0.407 | 0.235 | 0.145 | - |
| Speaker [58] | - | - | **0.204** | 0.494 |
| Gen | 0.353 | 0.194 | 0.120 | 0.415 |
| VC w/ Gen | 0.426 | 0.229 | 0.142 | 0.518 |
| VC w/ Gen+PG | **0.439** | **0.235** | 0.151 | **0.531** |
| | RefCOCO+ (Test B) | | | |
| | BLEU-1 | BLEU-2 | METEOR | CIDEr |
| MMI [30] | 0.324 | 0.167 | 0.133 | - |
| visdif [57] | 0.339 | 0.177 | 0.145 | - |
| Speaker [58] | - | - | **0.202** | 0.709 |
| Gen | 0.364 | 0.172 | 0.128 | 0.659 |
| VC w/ Gen | 0.391 | 0.197 | 0.146 | 0.731 |
| VC w/ Gen+PG | **0.404** | **0.209** | 0.154 | **0.742** |
| | RefCOCOg (Val*) | | | |
| | BLEU-1 | BLEU-2 | METEOR | CIDEr |
| MMI [30] | 0.428 | 0.263 | 0.144 | - |
| visdif [57] | 0.442 | 0.277 | 0.151 | - |
| Speaker [58] | - | - | **0.154** | 0.592 |
| Gen | 0.398 | 0.233 | 0.108 | 0.504 |
| VC w/ Gen | 0.456 | 0.281 | 0.139 | 0.625 |
| VC w/ Gen+PG | **0.467** | **0.287** | 0.146 | **0.630** |

TABLE 7
Human evaluation (Acc%) on referring expression generation. Note that [58] reports slightly higher accuracy using reranking mechanism. For fair comparison, we only report their performance without reranking.

| | RefCOCO | | RefCOCO+ | |
|---|---|---|---|---|
| | Test A | Test B | Test A | Test B |
| MMI [30] | 65.76 | 68.25 | 49.84 | 45.38 |
| Speaker [58] | 74.08 | 76.44 | 56.92 | 53.23 |
| Gen | 71.16 | 74.28 | 53.24 | 52.97 |
| VC w/ Gen | 74.52 | 77.18 | 56.04 | 56.26 |
| VC w/ Gen+PG | 74.39 | 77.56 | 56.35 | 56.48 |
| VC w/ Gen+PG (resnet) | **75.27** | **78.62** | **57.56** | **57.83** |

VGG feature. Fig. 10 presents some example generation results on three datasets. Our full model is shown to be able to generate concise and unambiguous description with important context information, such as location (*e.g.*, left, front), color (*e.g.*, red, blue), and related objects (*e.g.*, in red shirt, holding a dog). There are also cases that our model tends to fail. For the second example from the Test A split of RefCOCO+, our model succeeds to describe the audience in background using the clue "blurry" compared to batter in front, but fails to further distinguish the two audiences from color and location.

## 6 CONCLUSIONS

We focused on the task of grounding referring expressions in images and discussed that the key problem is how to model the complex context, which is not effectively resolved by the multiple instance learning framework used in prior works. Towards this challenge, we introduced the Variational Context framework, where the variational lower-bound can be interpreted by the reciprocity between the referent and context: given any of which can help to localize the other, and hence is expected to significantly reduce the context complexity in a principled way. The generation module is further included for semantic context modeling. The framework is implemented using cue-specific language-vision embedding network and policy gradient method that can be efficiently trained end-to-end. We validated the effectiveness of this reciprocity by promising supervised and unsupervised experiments on four benchmarks. We expect a future direction on one-stage visual grounding [4], [9], where the target object for the referring expression is directly localized without the region proposal generation stage for efficiency.

addition, our full model (VC w / Gen+PG) achieves the highest score on BLEU-1, BLEU-2 and CIDEr, while obtains slightly lower METEOR than Speaker.

Note that these metrics do not always reflect the ambiguity of generated description [57], since there are multiple possible expressions which can distinguish one object from others. Thus, we follow [57] and run a human evaluation on RefCOCO and RefCOCO+ to better evaluate the ambiguity of generated referring expression. Given the generated referring expression, the users were asked to click the referred object from the image, and the grounding accuracy was recorded for evaluation. The human evaluation results are shown in Table 7. Note that the comprehension module is only used for context estimation. The results show that generation module with estimated context has higher performance than vanilla generator, demonstrating that context modeling helps to generate unambiguous referring expression. In addition, the ResNet feature sightly improves the performance compared to the
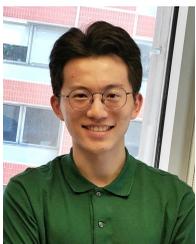
## REFERENCES

[1] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. Lawrence Zitnick, and D. Parikh. Vqa: Visual question answering. In *ICCV*, 2015.

[2] J. Ba, V. Mnih, and K. Kavukcuoglu. Multiple object recognition with visual attention. In *ICLR*, 2015.

[3] D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align and translate. 2015.

[4] X. Chen, L. Ma, J. Chen, Z. Jie, W. Liu, and J. Luo. Real-time referring expression comprehension by single-stage grounding network. *arXiv preprint arXiv:1812.03426*, 2018.

[5] K. Cho, B. Van Merriënboer, D. Bahdanau, and Y. Bengio. On the properties of neural machine translation: Encoder-decoder approaches. *arXiv preprint arXiv:1409.1259*, 2014.

[6] B. Dai, Y. Zhang, and D. Lin. Detecting visual relationships with deep relational networks. In *CVPR*, 2017.

GT  Ours  GT  Ours

**RefCOCO Test A**
- kid left / left kid
- bear in red / bear in red
- right bear / bear on right
- red bike / front bike
- far right bike / bike on right
- blue shirt / guy on bike

**RefCOCO Test B**
- pizza in front / the closest pizza
- pizza in the back / pizza in back
- donut in middle / donut in middle
- donut underneath / bottom donut on left
- donut on left
- top donut / top donut

**RefCOCO+ Test A**
- red shirt / man in red shirt
- man in white shirt / white shirt
- batter / batter
- man in red behind batter / blurry man in background
- white shirt background / blurry man in background

**RefCOCO+ Test B**
- closest bed / bed closest to us
- smaller bed / bed at 3 o'clock
- empty chair closest to us / chair closest to us
- totally empty chair / empty chair

**RefCOCOg Val\***
- a man in a blue coat standing in the snow / a man in blue jacket and black pants
- a pair of adult skis / the skis on the left
- the person wearing black who is holding a hot dog / a hand holding a dog
- gray pants of person standing to the back left of the hot dog / a person in the background
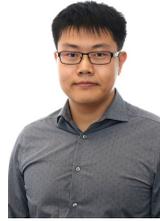
Fig. 10. Example generation results using our full model (VC w / Gen+PG) on three datasets. The ground-truth/generated expression is linked with the described referent using the same color.

[7] A. Das, S. Kottur, J. M. Moura, S. Lee, and D. Batra. Learning cooperative visual dialog agents with deep reinforcement learning. In *ICCV*, 2017.

[8] C. Deng, Q. Wu, Q. Wu, F. Hu, F. Lyu, and M. Tan. Visual grounding via accumulated attention. In *CVPR*, 2018.

[9] C. Deng, Q. Wu, G. Xu, Z. Yu, Y. Xu, K. Jia, and M. Tan. You only look & listen once: Towards fast and accurate visual grounding. *arXiv preprint arXiv:1902.04213*, 2019.

[10] T. G. Dietterich, R. H. Lathrop, and T. Lozano-Pérez. Solving the multiple instance problem with axis-parallel rectangles. *Artificial intelligence*, 1997.

[11] C. W. Fox and S. J. Roberts. A tutorial on variational bayesian inference. *Artificial intelligence review*, 2012.

[12] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In *ICAIS*, 2010.

[13] D. Golland, P. Liang, and D. Klein. A game-theoretic approach to generating spatial descriptions. In *EMNLP*, 2010.

[14] R. Hu, J. Andreas, M. Rohrbach, T. Darrell, and K. Saenko. Learning to reason: End-to-end module networks for visual question answering. *ICCV*, 2017.

[15] R. Hu, M. Rohrbach, J. Andreas, T. Darrell, and K. Saenko. Modeling relationships in referential expressions with compositional modular networks. In *CVPR*, 2017.

[16] R. Hu, H. Xu, M. Rohrbach, J. Feng, K. Saenko, and T. Darrell. Natural language object retrieval. In *CVPR*, 2016.

[17] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015.

[18] S. Kazemzadeh, V. Ordonez, M. Matten, and T. L. Berg. Referitgame: Referring to objects in photographs of natural scenes. In *EMNLP*, 2014.

[19] D. P. Kingma and M. Welling. Auto-encoding variational bayes. In *ICLR*, 2014.

[20] E. Krahmer and K. Van Deemter. Computational generation of referring expressions: A survey. *Computational Linguistics*, 38(1):173–218, 2012.

[21] Y. Li, N. Duan, B. Zhou, X. Chu, W. Ouyang, X. Wang, and M. Zhou. Visual question generation as dual task of visual question answering. In *CVPR*, 2018.

[22] Y. Li, W. Ouyang, and X. Wang. Vip-cnn: Visual phrase guided convolutional neural network. In *CVPR*, 2017.

[23] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014.

[24] J. Liu, L. Wang, and M.-H. Yang. Referring expression generation and comprehension via attributes. In *ICCV*, 2017.

[25] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg. Ssd: Single shot multibox detector. In *ECCV*, 2016.

[26] C. Lu, R. Krishna, M. Bernstein, and L. Fei-Fei. Visual relationship detection with language priors. In *ECCV*, 2016.

[27] J. Lu, J. Yang, D. Batra, and D. Parikh. Hierarchical question-image co-attention for visual question answering. In *NIPS*, 2016.

[28] R. Luo and G. Shakhnarovich. Comprehension-guided referring expressions. In *CVPR*, 2017.

[29] A. Makhzani, J. Shlens, N. Jaitly, and I. J. Goodfellow. Adversarial autoencoders. In *ICLR Workshop*, 2016.

[30] J. Mao, J. Huang, A. Toshev, O. Camburu, A. L. Yuille, and K. Murphy. Generation and comprehension of unambiguous object descriptions. In *CVPR*, 2016.

[31] M. Mitchell, K. van Deemter, and E. Reiter. Natural reference to objects in a visual domain. In *INLG*, 2010.

[32] M. Mitchell, K. Van Deemter, and E. Reiter. Generating expressions that refer to visible objects. In *NAACL*, 2013.

[33] V. K. Nagaraja, V. I. Morariu, and L. S. Davis. Modeling context between objects for referring expression understanding. In *ECCV*, 2016.

[34] J. Pennington, R. Socher, and C. Manning. Glove: Global vectors for word representation. In *EMNLP*, 2014.

[35] B. A. Plummer, A. Mallya, C. M. Cervantes, J. Hockenmaier, and S. Lazebnik. Phrase localization and visual relationship detection with comprehensive linguistic cues. In *ICCV*, 2017.

[36] B. A. Plummer, L. Wang, C. M. Cervantes, J. C. Caicedo, J. Hockenmaier, and S. Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In *ICCV*, 2015.

[37] J. Redmon and A. Farhadi. Yolo9000: better, faster, stronger. In *CVPR*, 2017.

[38] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NIPS*, 2015.

[39] A. Rohrbach, M. Rohrbach, R. Hu, T. Darrell, and B. Schiele. Grounding of textual phrases in images by reconstruction. In *ECCV*, 2016.

[40] M. Schuster and K. K. Paliwal. Bidirectional recurrent neural networks. *TSP*, 1997.

[41] S. Schuster, R. Krishna, A. Chang, L. Fei-Fei, and C. D. Manning. Generating semantically precise scene graphs from textual descriptions for improved image retrieval. In *Workshop on Vision and Language*, 2015.

[42] K. Sohn, H. Lee, and X. Yan. Learning structured output representation using deep conditional generative models. In *NIPS*, 2015.

[43] Q. Sun, B. Schiele, and M. Fritz. A domain based approach to social relation recognition. In *CVPR*, 2017.

[44] J. A. Thomas. *Meaning in interaction: An introduction to pragmatics*. Routledge, 2014.

[45] J. Thomason, J. Sinapov, and R. Mooney. Guiding interaction behaviors for multi-modal grounded language learning. In *Proceedings of the First Workshop on Language Grounding for Robotics*, pages 20–24, 2017.

[46] K. van Deemter, I. van der Sluis, and A. Gatt. Building a semantically transparent corpus for the generation of referring expressions. In *INLG*, 2006.

[47] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan. Show and tell: A neural image caption generator. In *CVPR*, 2015.

[48] L. Weaver and N. Tao. The optimal reward baseline for gradient-based reinforcement learning. In *UAI*, pages 538–545. Morgan Kaufmann Publishers Inc., 2001.

[49] Y. Wei, J. Feng, X. Liang, C. Ming-Ming, Y. Zhao, and S. Yan. Object region mining with adversarial erasing: A simple classification to semantic segmentation approach. In *CVPR*, 2017.

[50] R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.

[51] T. Winograd. Understanding natural language. *Cognitive psychology*, 3(1):1–191, 1972.

[52] F. Xiao, L. Sigal, and Y.-J. Lee. Weakly-supervised visual grounding of phrases with linguistic structures. In *CVPR*, 2017.

[53] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, and Y. Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *ICML*, 2015.

[54] T. Xue, J. Wu, K. Bouman, and B. Freeman. Visual dynamics: Probabilistic future frame synthesis via cross convolutional networks. In *NIPS*, 2016.

[55] X. Yan, J. Yang, K. Sohn, and H. Lee. Attribute2image: Conditional image generation from visual attributes. In *ECCV*, 2016.

[56] L. Yu, Z. Lin, X. Shen, J. Yang, X. Lu, M. Bansal, and T. L. Berg. Mattnet: Modular attention network for referring expression comprehension. In *CVPR*, 2018.

[57] L. Yu, P. Poirson, S. Yang, A. C. Berg, and T. L. Berg. Modeling context in referring expressions. In *ECCV*, 2016.

[58] L. Yu, H. Tan, M. Bansal, and T. L. Berg. A joint speaker-listener-reinforcer model for referring expressions. In *ICCV*, 2017.

[59] H. Zhang, Z. Kyaw, S.-F. Chang, and T.-S. Chua. Visual translation embedding network for visual relation detection. In *CVPR*, 2017.

[60] H. Zhang, Z. Kyaw, J. Yu, and S.-F. Chang. Ppr-fcn: Weakly supervised visual relation detection via parallel pairwise r-fcn. In *ICCV*, 2017.

[61] H. Zhang, Y. Niu, and S.-F. Chang. Grounding referring expressions in images by variational context. In *CVPR*, 2018.

[62] Z. Zhao, Q. Yang, D. Cai, X. He, and Y. Zhuang. Video question answering via hierarchical spatio-temporal attention networks. In *International Joint Conference on Artificial Intelligence (IJCAI)*, volume 2, 2017.

[63] B. Zhuang, Q. Wu, C. Shen, I. D. Reid, and A. van den Hengel. Parallel attention: A unified framework for visual object discovery through dialogs and queries. In *CVPR*, 2018.

[64] C. L. Zitnick and P. Dollár. Edge boxes: Locating object proposals from edges. In *ECCV*, 2014.

**Hanwang Zhang** is currently an Assistant Professor at Nanyang Technological University, Singapore. He was a research scientist at the Department of Computer Science, Columbia University, USA. He has received the B.Eng (Hons.) degree in computer science from Zhejiang University, Hangzhou, China, in 2009, and the Ph.D. degree in computer science from the National University of Singapore in 2014. His research interest includes computer vision, multimedia, and social media. Dr. Zhang is the recipient of the Best Demo runner-up award in ACM MM 2012, the Best Student Paper award in ACM MM 2013, and the Best Paper Honorable Mention in ACM SIGIR 2016, and TOMM best paper award 2018. He is also the winner of Best Ph.D. Thesis Award of School of Computing, National University of Singapore, 2014.

**Zhiwu Lu** received the M.S. degree in applied mathematics from Peking University in 2005, and the Ph.D. degree in computer science from the City University of Hong Kong in 2011. He is currently an Associate Professor with the School of Information, Renmin University of China. He received the Best Paper Award at CGI 2014 and the IBM SUR Award 2015. His research interests include machine learning, pattern recognition, and computer vision.

**Shih-Fu Chang** is the Senior Executive Vice Dean of engineering with Columbia University and a Professor in electrical engineering and computer science. His research interests include computer vision, machine learning, and multimedia information retrieval, with the goal to turn unstructured multimedia data into searchable information. He is a fellow of the American Association for the Advancement of Science and the ACM.

**Yulei Niu** received the B.E. degree in computer science from the Renmin University of China, Beijing, China, in 2015, where he is currently pursuing the Ph.D. degree in computer science. From 2017 to 2018, he visited the Digital Video and Multimedia Laboratory, Columbia University, as a Visiting Ph.D. Student, under the supervision of Prof. Shih-Fu Chang. His research interests include computer vision, multimedia, and machine learning.