# Learning to Assemble Neural Module Tree Networks for Visual Grounding

Daqing Liu[1]    Hanwang Zhang[2]    Feng Wu[1]    Zheng-Jun Zha[1*]

[1]University of Science and Technology of China, [2]Nanyang Technological University

liudq@mail.ustc.edu.cn, hanwangzhang@ntu.edu.sg, fengwu@ustc.edu.cn, zhazj@ustc.edu.cn

## Abstract

*Visual grounding, a task to ground (i.e., localize) natural language in images, essentially requires composite visual reasoning. However, existing methods over-simplify the composite nature of language into a monolithic sentence embedding or a coarse composition of subject-predicate-object triplet. In this paper, we propose to ground natural language in an intuitive, explainable, and composite fashion as it should be. In particular, we develop a novel modular network called Neural Module Tree network (NMTREE) that regularizes the visual grounding along the dependency parsing tree of the sentence, where each node is a neural module that calculates visual attention according to its linguistic feature, and the grounding score is accumulated in a bottom-up direction where as needed. NMTREE disentangles the visual grounding from the composite reasoning, allowing the former to only focus on primitive and easy-to-generalize patterns. To reduce the impact of parsing errors, we train the modules and their assembly end-to-end by using the Gumbel-Softmax approximation and its straight-through gradient estimator, accounting for the discrete nature of module assembly. Overall, the proposed NMTREE consistently outperforms the state-of-the-arts on several benchmarks. Qualitative results show explainable grounding score calculation in great detail.*

## 1. Introduction

Visual grounding (*a.k.a.*, referring expression comprehension) aims to localize a natural language description in an image. It is one of the core AI tasks for testing the machine comprehension of visual scene and language [18]. Perhaps the most fundamental and related grounding system for words is object detection [32] (or segmentation [8]): the image regions (or pixels) are classified to the corresponding word of the object class. Despite their diverse model architectures [22], their sole objective is to calculate a grounding score for a visual region and a word, measuring the semantic association between the two modalities.
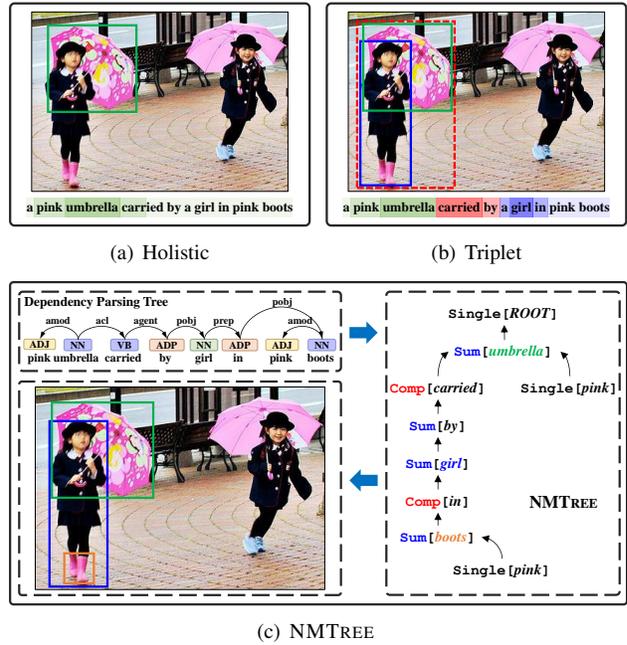


(a) Holistic    (b) Triplet



(c) NMTREE

Figure 1. Existing grounding models are generally (a) holistic or (b) coarsely composite. Words in gradient colors indicate word-level attentions. (c) The proposed NMTREE is based on dependency parsing tree and offers explainable grounding in great detail. Word color corresponds to image regions.

Thanks to the development of deep visual features [9] and language models [27], we can extend the grounding systems from fixed-size inventory of words to open-vocabulary [12] or even descriptive and relational phrases [41, 31].

However, grounding complex language sentences, *e.g.*, "a pink umbrella carried by a girl in pink boots", is far different from the above word or phrase cases. For example, given the image in Figure 1, for us humans, how to localize the "umbrella"? One may have the following reasoning process: 1) Identify the referent "umbrella", but there are two of them. 2) Use the contextual evidence "carried by a girl", but there are two girls. 3) By using more specific evidence "in pink boots", localize the "girl" in the last step. 4) Finally, by accumulating the above evidences, localize the target "umbrella".

1

Unfortunately, existing visual grounding methods generally rely on 1) a single monolithic score for the whole sentence [26, 38, 24, 39] (Figure 1(a)), or 2) a composite score for subject, predicate, and object phrases [13, 37] (Figure 1(b)). Though some of them adopt the word-level attention mechanism [25] to focus on the informative language parts, their reasoning is still coarse compared to the above human-level reasoning. More seriously, such coarse grounding scores are easily biased to learn certain vision-language patterns but not visual reasoning, *e.g.*, if most of the "umbrellas" are "carried by people" in the dataset, the score may not be responsive to other ones such as "people under umbrella stall". Not surprisingly, this problem has been repeatedly discovered in many end-to-end vision-language embedding frameworks used in other tasks such as VQA [15] and image captioning [23].

In this paper, we propose to exploit the Dependency Parsing Trees (DPTs) [3] that have already offered an off-the-shelf schema for the composite reasoning in visual grounding. Specifically, to empower the visual grounding ability by DPT, we propose a novel neural module network: Neural Module Tree (NMTREE) that provides explainable grounding scores in *great* detail. As illustrated in Figure 1(c), we transform a DPT into NMTREE by assembling three primitive module networks: Single for leaves and root, Sum and Comp for internal nodes (detailed in Section 3.3). Each module calculates a grounding score, which is accumulated in a bottom-up fashion, simulating the visual evidence gained so far. For example in Figure 1(c), Comp[$carried$] receives the scores gained by Sum[$by$] and then calculates a new score for the region composition, meaning "something is carried by the thing that is already grounded by the 'by' node". Thanks to the fixed reasoning schema, NMTREE disentangles the visual perception from the composite reasoning to alleviate the unnecessary vision-language bias [36], as the primitive modules receive consistent training signals with relatively simpler visual patterns and shorter language constitutions.

One maybe concerned by the potential brittleness caused by DPT parsing errors that impact the robustness of the module assembly, as discovered in most neural module networks applied in practice [11, 2]. We address this issue in three folds: 1) the assembly is simple. Except for Single that is fixed for leaves and root, only Sum and Comp are to be determined at run-time; 2) Sum is merely an Add operation that requires no visual grounding; 3) we adopt the recently proposed Gumbel-Softmax (GS) approximation [14] for the discrete assembly approximation. During training, the forward pass selects the two modules by GS sampler in a "hard" discrete fashion; the backward pass will update all possible decisions by using the straight-through gradient estimator in a "soft" robust way. By using the GS strategy, the entire NMTREE can be trained end-to-end without any

additional module layout annotations.

We validate the effectiveness of NMTREE on three challenging visual grounding benchmarks: RefCOCO [38], RefCOCO+ [38], and RefCOCOg [26]. NMTREE achieves new state-of-the-art performances on most of test splits and grounding tasks. Qualitative results and human evaluation indicate that NMTREE is transparent and explainable.

## 2. Related Work

Visual grounding is a task that requires a system to localize a region in an image while given a natural language expression. Different from object detection [32], the key for visual grounding is to utilize the linguistic information to distinguish the target from other objects, especially the objects of the same category.

To solve this problem, pioneering methods [26, 38, 24, 39] use the CNN-LSTM structure to localize the region that can generate the expression with maximum posteriori probability. Recently, joint embedding models [13, 37, 40] are widely used, they model the conditional probability and then localize the region with maximum probability conditioned on the expression. Our model belongs to the second category. However, compared with the previous works which neglect the rich linguistic structure, we step forward by taking structure information into account. Compared to [5] which relies on constituency parsing tree, our model applied dependency parsing tree with great parsing detail and the module assembly is learned end-to-end from scratch, while theirs is hand-crafted.

There are some works [13, 37] on using module networks in visual grounding task. However, they oversimplify the language structure and their modules are too coarse compared to ours. Fine-grained module networks are widely used in VQA [1, 2, 11]. However, they rely on additional annotations to learn a sentence-to-module layout parser, which is not available in general domains. Our module layout is trained from scratch by using the Gumbel-Softmax training strategy [14], which has shown empirically effective in recent works [4, 35, 29].

## 3. NMTREE Model

In this section, we first formulate the problem of visual grounding in Section 3.1. Then, by using the walk-through example illustrated in Figure 2, we introduce how to build NMTREE in Section 3.2 and how to calculate the grounding score using NMTREE in Section 3.3. Finally, we detail the Gumbel-Softmax training strategy in Section 3.4.

### 3.1. Problem Formulation

The visual grounding task can be reduced into a retrieval problem. Formally, given an image $\mathcal{I}$, we represent it by a set of Region of Interest (RoI) features $\mathcal{I} =$
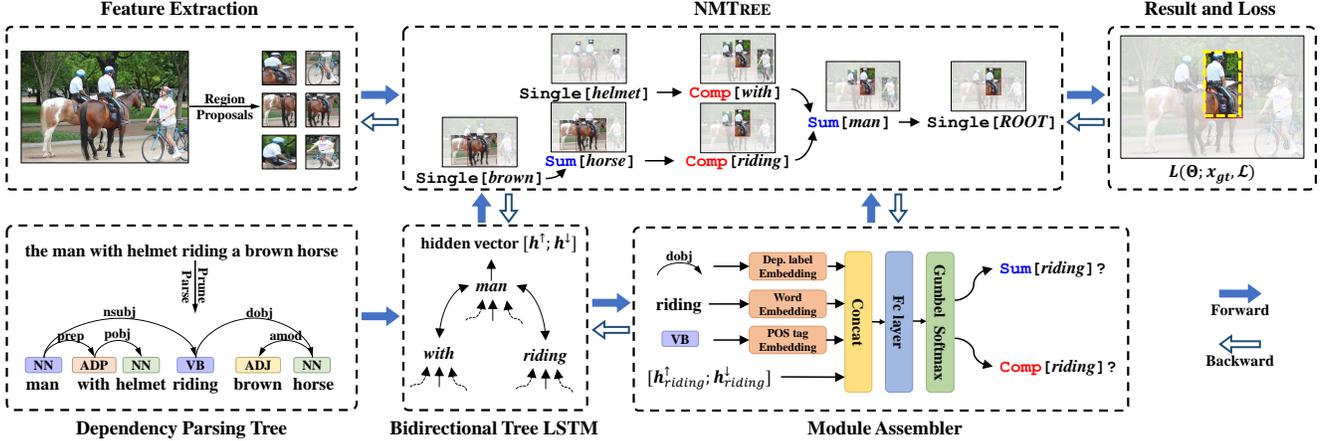
Figure 2. The overview of NMTREE for visual grounding. Given a natural language expression as input, we first transform it into NMTREE by Dependency Parsing Tree, Bidirectional Tree LSTM, and Module Assembler (Section 3.2). Then we ground along the tree in a bottom-up fashion (Section 3.3). The final result grounding score is the output score of the root node. We apply Gumbel-Softmax strategy to train our model (Section 3.4).

$\{\boldsymbol{x}_1, \boldsymbol{x}_2, \cdots, \boldsymbol{x}_K\}$, where $\boldsymbol{x}_i \in \mathbb{R}^{d_x}$ and $K$ is the number of regions. For a natural language phrase $\mathcal{L}$, we represent it by a word sequence $\mathcal{L} = \{w_1, w_2, \cdots, w_T\}$, where $T$ is the length of sentence. Then, the task is to retrieve the target region $\boldsymbol{x}^*$ by maximizing the grounding score $S(\boldsymbol{x}_i, \mathcal{L})$ between any region and the language:

$$\boldsymbol{x}^* = \arg\max_{\boldsymbol{x}_i \in \mathcal{I}} S(\boldsymbol{x}_i, \mathcal{L}). \tag{1}$$

Therefore, the key is to define a proper $S(\cdot)$ that distinguishes the target region from others by comprehending the language composition.

The pioneering grounding models [26, 38] are generally based on the holistic sentence-level language representation (Figure 1(a)): $S(\boldsymbol{x}_i, \mathcal{L}) := S_h(\boldsymbol{x}_i, \boldsymbol{y}_h)$, where $\boldsymbol{y}_h$ is a feature representation for the whole language expression and $S_h(\cdot)$ can be any similarity function between two vectors. More recently, a coarse composition [13] was proposed to represent the sentence as a (subject, relationship, object) triplet (Figure 1(b)). Thus, the score can be decomposed into a finer-grained composition: $S(\boldsymbol{x}_i, \mathcal{L}) := S_s(\boldsymbol{x}_i, \boldsymbol{y}_s) + S_r([\boldsymbol{x}_i, \boldsymbol{x}_o], \boldsymbol{y}_r) + S_o(\boldsymbol{x}_o, \boldsymbol{y}_o)$ where the subscripts $s$, $r$, and $o$ indicate the three linguistic roles: subject, relationship, and object, respectively; $\boldsymbol{x}_o$ is an estimated object region feature. However, these grounding scores oversimplify the composition of the language. For example, as shown in Figure 1(b), it is meaningful to decompose short sentences such as "umbrella carried by girl" into triplets, as it has a clear vision-language association for individual "girl", "umbrella", and their relationship; but it is problematic for longer sentences that are more general with clauses, e.g., even if the "girl in pink boots" is identified as the object, it is still coarse and difficult for grounding.

To this end, we propose to use the Dependency Parsing Tree (DPT) as a fine-grained language decomposition,

which empowers the grounding model to perform visual reasoning in great detail (Figure 1(c)):

$$S(\boldsymbol{x}_i, \mathcal{L}) := \sum\nolimits_t S_t(\boldsymbol{x}_i, \mathcal{L}_t), \tag{2}$$

where $t$ is a node in the tree, $S_t(\cdot)$ is a node-specific score function that calculates the similarity between a region and a node-specific language part $\mathcal{L}_t$. Intuitively, Eq. (2) is more human-like: accumulating the evidence (e.g., grounding score) while comprehending the language. Next, we will introduce how to implement Eq. (2).

## 3.2. Sentence to NMTREE

There are three steps to transform a sentence into the proposed NMTREE, as shown in the bottom three blocks of Figure 2. First, we parse the sentence into a DPT, where each word is a tree node. Then, we encode each word and its linguistic information into a hidden vector by a Bidirectional Tree LSTM. Finally, we assemble the neural modules to the tree according to node hidden vectors.

**Dependency Parsing Tree.** We adopt a dependency parser from Spacy toolbox[1]. As shown in Figure 2, it structures the language into a tree, where each node is a word with its part-of-speech (POS) tag and dependency relation label of the directed edge from it to another, e.g., "riding" is VB (verb) and its nsubj (nominal subject) is "man" as NN (noun). DPT offers an in-depth comprehension of a sentence and its tree structure offers a reasoning path for visual grounding. Note that there are always unnecessary syntax elements parsed from a free-form sentence such as determiners, symbols, and punctuation. We remove these nodes and edges to reduce the computational overhead without hurting the performance.

---

[1]Spacy2: https://spacy.io/

**Bidirectional Tree LSTM.** Once the DPT is obtained, we encode each node into a hidden vector by a bidirectional tree-structured LSTM [34]. This bidirectional (*i.e.*, bottom-up and top-down) propagation makes each node being aware of the information both from its children and parent. This is particularly crucial for capturing the context in a sentence. For each node $t$, we embed the word $w_t$, POS tag $p_t$, and dependency relation label $d_t$ into a concatenated embedding vector as:

$$\boldsymbol{e}_t = [\mathbf{E}_w \Pi_{w_t}, \mathbf{E}_p \Pi_{p_t}, \mathbf{E}_d \Pi_{d_t}], \quad (3)$$

where $\mathbf{E}_w$, $\mathbf{E}_p$, and $\mathbf{E}_d$ are trainable embedding matrices, $\Pi_{w_t}$, $\Pi_{p_t}$, and $\Pi_{d_t}$ are one-hot encodings, for word, POS tag, and dependency relation label, respectively.

Our tree LSTM implementation [2] is based on the Child-Sum Tree LSTM [34]. Taking the bottom-up direction for example, a node $t$ receives the LSTM states from its children node set $\mathcal{C}_t$ and its embedding vector $\boldsymbol{e}_t$ as input to update the state:

$$\boldsymbol{c}_t^\uparrow, \boldsymbol{h}_t^\uparrow = \text{TreeLSTM}(\boldsymbol{e}_t, \{\boldsymbol{c}_{tj}^\uparrow\}, \{\boldsymbol{h}_{tj}^\uparrow\}), \quad j \in \mathcal{C}_t, \quad (4)$$

where $\boldsymbol{c}_{tj}^\uparrow$, $\boldsymbol{h}_{tj}^\uparrow$ denote the cell and hidden vectors of the $j$-th child of node $t$. By applying the TreeLSTM in two directions, we can obtain the final node hidden vector $\boldsymbol{h}_t$ as:

$$\boldsymbol{h}_t = [\boldsymbol{h}_t^\uparrow; \boldsymbol{h}_t^\downarrow], \quad (5)$$

where $\boldsymbol{h}_t^\uparrow$, $\boldsymbol{h}_t^\downarrow \in \mathbb{R}^{d_h}$ denote the hidden vectors encoded in the bottom-up and top-down directions, respectively. We initialize all leaf nodes with zero hidden and cell states. The bottom-up and top-down Tree LSTMs have their independent trainable parameters.

**Module Assembler.** Given the node representation $\boldsymbol{e}_t$ and the above obtained node hidden vector $\boldsymbol{h}_t$, we can feed them into a module assembler to determine which module should be assembled to node $t$. As we will detail in Section 3.3, we have three modules, *i.e.*, Single, Sum, and Comp. Since the Single is always assembled on leaves and the root, the assembler only need to choose between Sum and Comp as:

$$\text{Sum or Comp} \leftarrow \arg\max \text{softmax} \left(\text{fc}([\boldsymbol{e}_t, \boldsymbol{h}_t])\right), \quad (6)$$

where fc is a fully connected layer that maps the input features into a 2-d values, indicating the relative scores for Sum and Comp, respectively. Due to the discrete and non-differentiable nature of arg max, we use the Gumbel-Softmax [14] strategy for training (Section 3.4).

It is worth noting that the assembler is not purely linguistic even though Eq. (6) is based on DPT node features. In fact, thanks to the back-propagation training algorithm, visual cues will be eventually incorporated into the parameters of Eq. (6). Figure 3 illustrates which type of words is



(a) Sum Module Node     (b) Comp Module Node

Figure 3. Word cloud visualizations of what nodes are likely to be assembled as Sum or Comp. We can find that the Sum module nodes are likely to be visible concepts while the Comp module nodes are likely to be relationship concepts.

likely to be assembled by each module. We can find that the Sum module has more visible words (*e.g.*, adjectives and nouns), and the Comp module has more words describing relations (*e.g.*, verbs and prepositions). This reveals the explainable potential of NMTREE. Finally, by the above three steps, we get the NMTREE that each node assembled. Next, we will elaborate the three types of modules.

### 3.3. NMTREE **Modules**

Given the above assembled NMTREE, we can implement the tree grounding score proposed in Eq. (2) by accumulating the scores in a bottom-up fashion. There are three types of modules used in NMTREE, *i.e.*, Single, Sum and Comp. Each module at node $t$ updates the grounding score $\boldsymbol{s}_t = [s_t^1, \cdots, s_t^K]$ for all the $K$ regions in the image $\mathcal{I}$ and outputs to its parent. In the following, we will first introduce language representation and common functions used in the modules, and then detail each module.

**Language Representation.** For node $t$, we have two language representations: $\boldsymbol{y}_t^s$ is used to associate with a single visual feature and $\boldsymbol{y}_t^p$ is used to associate with a pairwise visual feature. We denote the node set of node $t$ as $\mathcal{N}_t$, which contains itself and all nodes rooted from $t$. Therefore, the language representation can be calculated by the weighted sum of node embedding vectors from $\mathcal{N}_t$:

$$\boldsymbol{y}_t^s = \sum_{i \in \mathcal{N}_t} \alpha_i^s \boldsymbol{e}_i, \quad \boldsymbol{y}_t^p = \sum_{i \in \mathcal{N}_t} \alpha_i^p \boldsymbol{e}_i, \quad (7)$$

where $\alpha$ are the node-level attention weights that calculated from the corresponding node hidden vectors: $\alpha_i = \text{softmax}(\text{fc}(\boldsymbol{h}_i))$. Note that $\alpha_i^s$ and $\alpha_i^p$ have independent fc parameters. It is worth noting that these weighted average word embeddings of the node set reduce the negative impact caused by DPT parsing errors [13].

**Score Functions.** There are two types of score functions used in our modules, denoted by the single score function $S_s$ and pairwise score function $S_p$, where $S_s$ measures the similarity between a single region $\boldsymbol{x}$ and a language representation $\boldsymbol{y}$, and $S_p$ indicates how likely pair-wise regions

---

[2] For space reasons, we leave the details in supplementary material.

4

match with relationships. Formally we define them as:

$$S_s(\boldsymbol{x}, \boldsymbol{y}) = \text{fc}(\text{L2norm}(\text{fc}(\boldsymbol{x}) \odot \boldsymbol{y})), \qquad (8)$$

$$S_p(\boldsymbol{x}_1, \boldsymbol{x}_2, \boldsymbol{y}) = \text{fc}(\text{L2norm}(\text{fc}([\boldsymbol{x}_1; \boldsymbol{x}_2]) \odot \boldsymbol{y})), \qquad (9)$$

where $[;]$ is a concatenation operation, $\odot$ is element-wise multiplication, and L2norm is used to normalize vectors.

**`Single` Module**. It is assembled at leaves and the root. Its job is to 1) calculate a single score for each region and the current language feature by Eq. (8), 2) add this new score to the scores collected from children, and then 3) pass the sum to its parent:

$$
\begin{aligned}
\textbf{Input:} \quad & \{\boldsymbol{s}_{tj}\}, \quad j \in \mathcal{C}_t \\
\textbf{Output:} \quad & s_t^i \leftarrow S_s(\boldsymbol{x}_i, \boldsymbol{y}_t^s) + \sum_j s_{tj}^i, i \in [1, K]
\end{aligned}
\qquad (10)
$$

Note that for leaves, $\mathcal{C}_t = \phi$ as they have no children. As illustrated in Figure 2, its design motivation is to initiate the bottom-up grounding process by the most elementary words and finalize the grounding by passing the accumulated scores to ROOT.

**`Sum` Module**. It plays a transitional role during the reasoning process. It simply sums up the scores passed from its children and then passes the sum to its parent:

$$
\begin{aligned}
\textbf{Input:} \quad & \{\boldsymbol{s}_{tj}\}, \quad j \in \mathcal{C}_t \\
\textbf{Output:} \quad & \boldsymbol{s}_t \leftarrow \sum_j \boldsymbol{s}_{tj}
\end{aligned}
\qquad (11)
$$

Note that this module has no parameters hence it significantly reduces the complexity of our model. As illustrated in Figure 2, intuitively, it transits the easy-to-localize words (cf. Figure 3(a)) such as "horse" and "man" to help the subsequent composite grounding.

**`Comp` Module.** This is the core module for composite visual reasoning. As shown in Figure 3(b), it is likely to be the relationship that connects two language constitutions. It first computes an "average region" visual feature that is grounded by the single scores:

$$\beta_i = \text{softmax}\Big(S_s(\boldsymbol{x}_i, \boldsymbol{y}_t^s) + \sum_j s_{tj}^i\Big), \quad \bar{\boldsymbol{x}} = \sum_i \beta_i \boldsymbol{x}_i. \quad (12)$$

In particular, $\bar{\boldsymbol{x}}$ can be considered as the contextual region [42] that supports the target region score, *e.g.*, "what is riding the horse" in Figure 2. Therefore, this module outputs the target region score to its parent:

$$
\begin{aligned}
\textbf{Input:} \quad & \{\boldsymbol{s}_{tj}\}, \quad j \in \mathcal{C}_t \\
\textbf{Output:} \quad & s_t^i \leftarrow S_p(\boldsymbol{x}_i, \bar{\boldsymbol{x}}, \boldsymbol{y}_t^p).
\end{aligned}
\qquad (13)
$$

Recall that $\mathbf{y}_t^p$ is pairwise language feature that represents the relationship words.

By reasoning along the assembled NMTREE in bottom up fashion, we can obtain the overall accumulated grounding score in Eq. (2) at tree root. Moreover, thanks to the



Figure 4. A very long sentence example that illustrates the explainability of the neural modules in NMTREE. Black words:`Single`, blue words: `Sum`, red words: `Comp`.

score output at each node, NMTREE is transparent as the scores can be visualized as attention maps to investigate the grounding process. Figure 4 illustrates an extreme example with a very long expression with 22 tokens. However, by using the neural modules in NMTREE , it still works well and reasons with explainable intermediate process. Next, we will discuss how to train NMTREE.

### 3.4. NMTREE **Training**

In contrast to previous neural module networks [1, 11], NMTREE does not require any additional annotations and is end-to-end trainable. Suppose $\boldsymbol{x}_{gt}$ is the ground-truth region, the objective is to minimize the cross-entropy loss:

$$L(\Theta; \boldsymbol{x}_{gt}, \mathcal{L}) = -\log \text{softmax}(S(\boldsymbol{x}_{gt}, \mathcal{L}; \Theta)), \qquad (14)$$

where $\Theta$ is the trainable parameter set and softmax is across all $K$ regions in an image.

Recall that the assembling process in Eq. (6) is discrete and blocks the end-to-end training. Therefore, we utilize the Gumbel-Softmax strategy [7] that is shown effective in recent works [4, 35] on architecture search. For more details, please refer to their papers. Here, we only introduce how to apply the Gumbel-Softmax for NMTREE training.

**Forward**. We add Gumbel distribution as a noise into the relative scores (*i.e.* $\text{fc}([\boldsymbol{e}_t, \boldsymbol{h}_t])$) of each module. It introduces stochasticity for the module assembling exploration. Specifically, we parameterize the assembler decision as a 2-d one-hot vector $\boldsymbol{z}$, where the index of non-zero entry indicates the decision:

$$\boldsymbol{z} = \texttt{one\_hot}(\arg\max(\log(\text{fc}([\boldsymbol{e}_t, \boldsymbol{h}_t])) + G)), \quad (15)$$

where $G$ is the noise drawn from i.i.d. Gumbel(0, 1)[3]. Note that, in inference phrase, the $G$ will be discarded.

**Backward.** We take a continuous approximation that relaxes $\boldsymbol{z}$ to $\widetilde{\boldsymbol{z}}$ by replacing argmax with softmax, formally:

$$\widetilde{\boldsymbol{z}} = \text{softmax}((\log(\text{fc}([\boldsymbol{e}_t, \boldsymbol{h}_t])) + G)/\tau), \qquad (16)$$

---

[3]The Gumbel (0, 1) distribution is sampled by $G = -\log(-\log(U))$ where $U \sim \text{Uniform}(0, 1)$.

where $G$ is the same sample drawn in the forward pass (*i.e.*, we reuse the noise samples). $\tau$ is a temperature parameter that the softmax function approaches to argmax while $\tau \to 0$ and approaches to uniform while $\tau \to \infty$. Although there are discrepancies between the forward and backward pass, we empirically observe that the Gumbel-Softmax strategy performs well in our experiments.

## 4. Experiments

### 4.1. Datasets

We conducted our experiments on three datasets that are collected from MS-COCO [20] images. **RefCOCO** [38] contains 142,210 referring expressions for 19,994 images. An interactive game [16] is used during the expression collection. All expression-referent pairs are split into train, validation, testA, and testB. TestA contains the images with multiple people and testB contains the images with multiple objects. **RefCOCO+** [38] contains 141,564 referring expressions for 49,856 objects in 19,992 images. It is collected with the same interactive game as RefCOCO and is split into train, validation, testA, and testB, respectively. The difference from RefCOCO is that RefCOCO+ only allows expression described by appearance but no locations. **RefCOCOg** [26] contains 95,010 referring expressions for 49,822 objects in 25,799 images. It is collected in a non-interactive way and contains longer expressions described by both appearance and locations. It has two types of data partitions. The first partition [26] divides dataset into train and validation (val$^*$) sets. The second partition [28] divides images into train, validation (val) and test sets.

### 4.2. Implementation Details and Metrics

**Language Settings.** We built specific vocabularies for the three datasets with words, POS tags, and dependency labels appeared more than once in datasets. Note that to obtain accurate parsing results, we did not trim the length of expressions. We used pre-trained GloVe [30] to initialize word vectors. For dependency label vectors and POS tag vectors, we trained them from scratch with random initialization. We set the embedding sizes to 300, 50, 50 for words, POS tags, and dependency labels, respectively.

**Visual Representations.** To represent RoI features of an image, we concatenated object features and location features extracted from MAttNet [37], which is based on Faster RCNN [32] with ResNet-101 [9] as the backbone and trained with attribute heads. We employed Mask RCNN [8] for object segmentation. The visual feature dimension $d_x$ was set to 3,072. For fair comparison, we also used VGG-16 [33] as the backbone and $d_x$ was set to 5,120.

**Parameter Settings.** We optimized our model with Adam optimizer [17] up to 40 epochs. The learning rate was initialized to 1e-3 and shrunk by 0.9 every 10 epochs. We set

128 images to the mini-batch size. The LSTM hidden size $d_h$ was set to 1,024, the hidden size of the attention in language representation was set to 1,024. The temperature $\tau$ of Gumbel-Softmax [14] was set to 1.0.

**Evaluation Metrics.** For detection task, we calculated the Intersection-over-Union (IoU) between the detected bounding box and the ground-truth one, and treated the one with IoU at least 0.5 as correct. We used the Top-1 accuracy as the metric, which is the fraction of the correctly grounded test expressions. For segmentation task, we used Pr@0.5 (the percentage of expressions where IoU at least 0.5) and overall IoU as metrics.

### 4.3. Ablation Studies

**Settings**. We conducted extensive ablation studies to reveal the internal mechanism of NMTREE. The ablations and their motivations are detailed as follows. **Chain**: it ignores the structure information of the language. Specifically, we represent a natural language expression as the weighted average of each word embedding, where the weights are calculated by soft attention on bi-LSTM hidden vectors of each word. The final grounding score is calculated by single score function between each region and the language representation. **NMTREE w/o `Comp`**: it is the NMTREE without the `Comp` module, forcing all internal nodes as `Sum` module. **NMTREE w/o `Sum`**: it is the NMTREE without the `Sum` module, forcing all internal nodes as `Comp` module. **NMTREE w/ Rule**: it assembles modules by a hand-crafted rule. Instead of deciding which module should be assembled to each node by computing the relative score, we designed a fixed linguistic rule to make a discrete and non-trainable decisions. The rule is: set the internal nodes whose dependency relation label is 'acl' (*i.e.*, adjectival clause) or 'prep' (*i.e.*, prepositional modifier) as `Comp` module, and the others as `Sum` module.

**Results.** Table 1 shows the grounding accuracies of the ablation methods on the three benchmarks. We can have the following observations: 1) On all datasets, NMTREE outperforms Chain even if we removed one module or used the hand-crafted rule. This is because the tree structure contains more linguistic information and more suitable for reasoning. Meanwhile, it also demonstrates that our proposed fine-grained composition is better than the holistic Chain. 2) When we removed one module, *i.e.*, NMTREE w/o `Comp` and NMTREE w/o `Sum`, they are worse than the full NMTREE. It demonstrates the necessity of the `Sum` and `Comp`. Note that removing any modules will also hurt the explainability of models. 3) NMTREE w/o `Comp` and NMTREE w/o `Sum` are comparable but NMTREE w/o `Sum` is slightly better. This is because the `Comp` module is more complex and thus resulting in overfitting. 4) NMTREE outperforms NMTREE w/ Rule. It demonstrates that NMTREE can automatically find which nodes need composite reason-

| | RefCOCO | | | RefCOCO+ | | | RefCOCOg | | RefCOCO(det) | | | RefCOCO+(det) | | | RefCOCOg(det) | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | val | testA | testB | val | testA | testB | val | test | val | testA | testB | val | testA | testB | val | test |
| Chain | 82.43 | 82.21 | 82.16 | 68.27 | 70.83 | 62.41 | 73.84 | 74.15 | 74.81 | 79.19 | 68.34 | 63.08 | 68.84 | 53.53 | 61.72 | 61.95 |
| NMTREE w/o Comp | 83.65 | 83.59 | 83.04 | 70.76 | 73.07 | 65.19 | 75.98 | 76.20 | 75.10 | 79.38 | 68.60 | 64.85 | 70.43 | 55.00 | 63.07 | 63.40 |
| NMTREE w/o Sum | 83.79 | 83.81 | 83.67 | 70.83 | 73.72 | 65.83 | 76.11 | 76.09 | 75.49 | 79.84 | 69.11 | 65.29 | 70.85 | 55.99 | 63.60 | 64.06 |
| NMTREE w/ Rule | 84.46 | 84.59 | 84.26 | 71.48 | 74.76 | 66.95 | 77.82 | 77.70 | 75.51 | 80.61 | 69.23 | 65.23 | 70.94 | 56.96 | 64.69 | 65.53 |
| NMTREE | **85.65** | **85.63** | **85.08** | **72.84** | **75.74** | **67.62** | **78.57** | **78.21** | **76.41** | **81.21** | **70.09** | **66.46** | **72.02** | **57.52** | **65.87** | **66.44** |

Table 1. Top-1 Accuracy% of ablation models on the three datasets.

| | RefCOCO | | | RefCOCO+ | | | RefCOCOg | | | RefCOCO (det) | | | RefCOCO+ (det) | | | RefCOCOg (det) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | val | testA | testB | val | testA | testB | val* | val | test | val | testA | testB | val | testA | testB | val* | val | test |
| MMI [26] | - | 63.15 | 64.21 | - | 48.73 | 42.13 | 62.14 | - | - | - | 64.90 | 54.51 | - | 54.03 | 42.81 | 45.85 | - | - |
| Attribute [21] | - | 78.85 | 78.07 | - | 61.47 | 57.22 | 69.83 | - | - | - | 72.08 | 57.29 | - | 57.97 | 46.20 | 52.35 | - | - |
| Listener† [39] | 78.36 | 77.97 | 79.86 | 61.33 | 63.10 | 58.19 | 72.02 | 71.32 | 71.72 | 68.95 | 72.95 | 62.98 | 54.89 | 59.61 | 48.44 | 58.32 | 59.33 | 59.21 |
| NegBag [28] | 76.90 | 75.60 | 78.80 | - | - | - | - | - | 68.40 | 57.30 | 58.60 | 56.40 | - | - | - | 39.50 | - | 49.50 |
| CMN [13] | - | 75.94 | 79.57 | - | 59.29 | 59.34 | 69.30 | - | - | - | 71.03 | 65.77 | - | 54.32 | 47.76 | 57.47 | - | - |
| VC [42] | - | 78.98 | 82.39 | - | 62.56 | 62.90 | 73.98 | - | - | - | 73.33 | 67.44 | - | 58.40 | 53.18 | 62.30 | - | - |
| AccumAttn [6] | 81.27 | 81.17 | 80.01 | 65.56 | 68.76 | 60.63 | 73.18 | - | - | - | - | - | - | - | - | - | - | - |
| MAttN‡ [37] | 85.65 | 85.26 | 84.57 | 71.01 | 75.13 | 66.17 | - | 78.10 | 78.12 | 76.40 | 80.43 | 69.28 | 64.93 | 70.26 | 56.00 | - | **66.67** | **67.01** |
| GroundNet [5] | - | - | - | - | - | - | 68.90 | - | - | - | - | - | - | - | - | - | - | - |
| parser+CMN [13] | - | - | - | - | - | - | 53.50 | - | - | - | - | - | - | - | - | - | - | - |
| parser+MAttN‡ [37] | 80.20 | 79.10 | 81.22 | 66.08 | 68.30 | 62.94 | - | 73.82 | 73.72 | - | - | - | - | - | - | - | - | - |
| NMTREE | 80.39 | 78.86 | 81.90 | 63.31 | 63.59 | 63.04 | 73.71 | 73.39 | 72.29 | 71.65 | 74.81 | 67.34 | 58.00 | 61.09 | 53.45 | 61.20 | 61.01 | 61.46 |
| NMTREE ‡ | **85.65** | **85.63** | **85.08** | **72.84** | **75.74** | **67.62** | **78.03** | **78.57** | **78.21** | **76.41** | **81.21** | **70.09** | **66.46** | **72.02** | **57.52** | **64.62** | 65.87 | 66.44 |

Table 2. Top-1 Accuracy% of various grounding models on the three datasets. For fair comparison, we use † to indicate that this model uses res101 features for detected experiments. ‡ indicates that this model uses res101 features for both ground-truth and detected experiments. None-superscript indicates that this model uses vgg16 features. Parser+ indicates that the model used an external parser.

| | | RefCOCO | | | RefCOCO+ | | | RefCOCOg | |
|---|---|---|---|---|---|---|---|---|---|
| | | val | testA | testB | val | testA | testB | val | test |
| Pr@0.5 | MAttNet [37] | **75.16** | 79.55 | 68.87 | 64.11 | 70.12 | 54.82 | **64.48** | **65.60** |
| | Chain | 73.36 | 77.55 | 67.30 | 61.60 | 67.15 | 52.24 | 59.64 | 60.29 |
| | NMTREE | 74.71 | **79.71** | **68.93** | **65.06** | **70.24** | **56.15** | 63.77 | 64.63 |
| IoU | MAttNet [37] | 56.51 | 62.37 | 51.70 | 46.67 | 52.39 | 40.08 | **47.64** | **48.61** |
| | Chain | 55.29 | 60.99 | 51.36 | 44.74 | 49.83 | 38.50 | 42.55 | 43.99 |
| | NMTREE | **56.59** | **63.02** | **52.06** | 47.40 | **53.01** | 41.56 | 46.59 | 47.88 |

Table 3. Segmentation performance(%) on the three datasets comparing with state-of-the-arts.

ing (as Comp) or not (as Sum). Further, it also implies that our NMTREE is more suitable for visual grounding task as our assembler is aware of visual cues by the Gumbel-Softmax training strategy.

### 4.4. Comparison with State-of-the-Arts

**Settings**. We compared NMTREE with other state-of-the-art visual grounding models published in recent years. According to whether the model requires language composition, we group those methods into: 1) Generation based methods which select the region with the maximum generation probability: **MMI** [26], **Attribute** [21], and **Listener** [39]. 2) Holistic language based methods: **NegBag** [28]. 3) Language composition based methods: **CMN** [13], **VC** [42], **AccumAttn** [6], and **MAttN** [37]. 4) Composition methods with external parser: **GroundNet** [5], **parser+CMN**, and **parser+MAttN**. NMTREE belongs to the fourth category, but its language composition is more fine-grained than others. We compared with them on three different settings: ground-truth regions, detected regions, and segmentation masks.

**Results**. From Table 2 and Table 3, we can find that: 1) the

triplet composition models mostly outperform holistic models. This is because taking the advantage of linguistics information by decomposing sentences, even coarse-grained, is helpful in visual grounding. 2) Our model outperforms most triplet models with the help of fine-grained composite reasoning. 3) The parser-based methods are fragile to parser errors, leading to performance decline. However, our model is more robust because of the dynamic assembly and end-to-end train strategy. Although some of the performance gains are marginal, one should notice that it seems NMTREE balances the well-known trade-off between performance and explainability [10]. As we will discuss in the following, we achieve the explainability without hurting the accuracy.

### 4.5. Qualitative Analysis

In this section, we would like to investigate the internal reasoning steps of our model by qualitative results[4]. In Figure 5, we visualize the tree structures, the module assembly, the attention map at each intermediate step, and the final results. In Figure 6, we visualize the reasoning process inside Comp modules. With these qualitative visualizations, we can have the following observations: 1) The visual concept words usually are assembled by Sum module while the relationship concept words are usually assembled by Comp module. 2) The attention maps of non-visual leaf nodes, *e.g.*, 'directly' in 5(d), are usually scattered, while visual ones, *e.g.*, 'girl' in 5(d), are usually concentrated. 3) Comp modules are aware of relationships, *i.e.*, it can move the at-

---

[4]Since our work focuses on complex language cases, we mainly conducted qualitative experiments on RefCOCOg. More qualitative results are given in supplementary material.
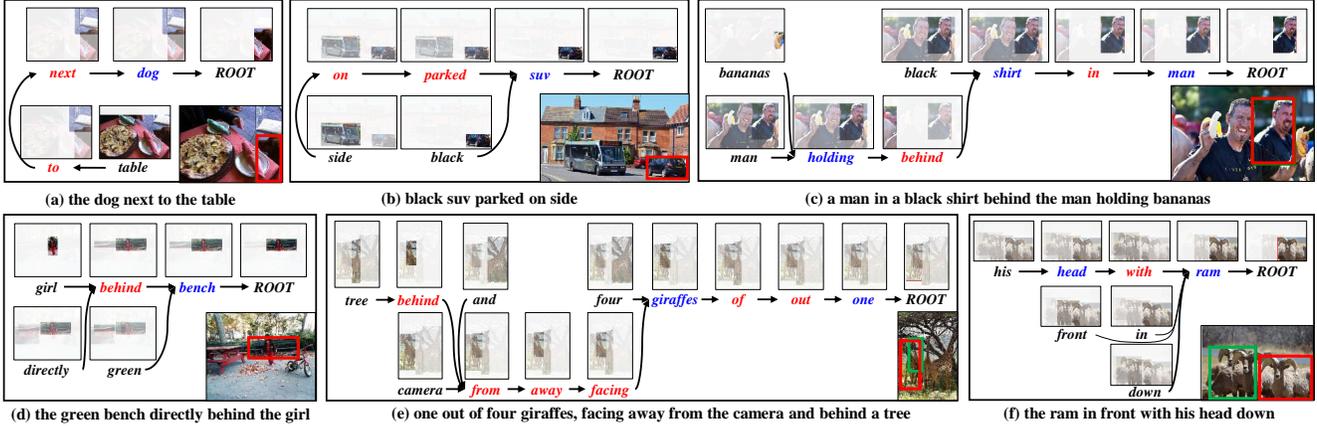
(a) the dog next to the table  (b) black suv parked on side  (c) a man in a black shirt behind the man holding bananas

(d) the green bench directly behind the girl  (e) one out of four giraffes, facing away from the camera and behind a tree  (f) the ram in front with his head down

Figure 5. Qualitative results on RefCOCOg. Words in different colors indicate corresponding modules: black for `Single`, red for `Comp`, and blue for `Sum`. The bottom right corner is the original image with a green bounding box as ground-truth and a red bounding box as our result. We further give two failure examples (e) and (f) for comparison, and our model consistently provides explainable reasoning process.



man holding blue surfboard  |  man in colored Beanie riding green bike

blue bus behind police man  |  zebra that is to left of rightmost zebra

Figure 6. The compositional reasoning inside `Comp`. Each example contains the original image (left), the contextual attention map of $\bar{x}$ in Eq. (12) (middle) and the output attention map (right). We represent partial tree structure by colors: red for the current node, blue for children and green for parent.



Figure 7. The percentage of each choice. The average scores of NMTREE and AccumAttn [6] are **2.96** and **2.28**. The results indicate that our model is more explainable to humans.

tention from the supporting objects to the target objects, as shown in Figure 6. 4) Along the tree, attention maps become more sharp, indicating the confidence of our model become stronger.

All the above observations suggest that our NMTREE can reason along the tree and provide rich cues to support the final results. These reasoning patterns and supporting cues imply that our model is explainable. Therefore, to further investigate the explainability of our model, we conducted a human evaluation to measure whether the inter-

nal reasoning process is reasonable. Since the state-of-the-art model MAttNet [37] does not contain internal reasoning process but only sums up three pre-defined module scores which directly point to the desired object, we compared with AccumAttn [6] for it performs multi-step sequential reasoning and has image/textual attention at each time step. We first presented 60 examples with internal steps of each model to 6 human evaluators, and asked them to judge how clear that the model was doing at each step. Then each evaluator rated each example on 4-point Likert scale [19] (unclear, slightly clear, mostly clear, clear) corresponding to scores of 1, 2, 3, and 4. The percentage of each choice and average scores are shown in Figure 7. We can find that our model outperforms AccumAttn [6] and is often rated as "clear". It indicates that the internal reasoning process of our model can be more clearly understood by humans.

## 5. Conclusion

In this paper, we proposed Neural Module Tree Networks (NMTREE), a novel end-to-end model that localizes the target region by accumulating the grounding confidence score along the dependency parsing tree of a natural language sentence. NMTREE consists of three simple neural modules, whose assembly is trained without additional annotations. Compared with previous visual grounding methods, our model performs a more fine-grained and explainable language composite reasoning with superior performance, demonstrated by extensive experiments on three benchmarks.

# References

[1] Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Dan Klein. Neural module networks. In *CVPR*, 2016. 2, 5

[2] Qingxing Cao, Xiaodan Liang, Bailing Li, Guanbin Li, and Liang Lin. Visual question reasoning on general dependency tree. In *CVPR*, 2018. 2

[3] Danqi Chen and Christopher Manning. A fast and accurate dependency parser using neural networks. In *EMNLP*, 2014. 2

[4] Jihun Choi, Kang Min Yoo, and Sang-goo Lee. Learning to compose task-specific tree structures. In *AAAI*, 2018. 2, 5

[5] Volkan Cirik, Taylor Berg-Kirkpatrick, and Louis-Philippe Morency. Using syntax to ground referring expressions in natural images. In *AAAI*, 2018. 2, 7

[6] Chaorui Deng, Qi Wu, Qingyao Wu, Fuyuan Hu, Fan Lyu, and Mingkui Tan. Visual grounding via accumulated attention. In *CVPR*, 2018. 7, 8

[7] Emil Julius Gumbel. Statistical theory of extreme values and some practical applications. *NBS Applied Mathematics Series*, 33, 1954. 5

[8] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *ICCV*, 2017. 1, 6

[9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 1, 6

[10] Ronghang Hu, Jacob Andreas, Trevor Darrell, and Kate Saenko. Explainable neural computation via stack neural module networks. In *ECCV*, 2018. 7

[11] Ronghang Hu, Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Kate Saenko. Learning to reason: End-to-end module networks for visual question answering. In *ICCV*, 2017. 2, 5

[12] Ronghang Hu, Piotr Dollár, Kaiming He, Trevor Darrell, and Ross Girshick. Learning to segment every thing. In *CVPR*, 2018. 1

[13] Ronghang Hu, Marcus Rohrbach, Jacob Andreas, Trevor Darrell, and Kate Saenko. Modeling relationships in referential expressions with compositional modular networks. In *CVPR*, 2017. 2, 3, 4, 7

[14] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. In *ICLR*, 2017. 2, 4, 6

[15] Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *CVPR*, 2017. 2

[16] Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. Referitgame: Referring to objects in photographs of natural scenes. In *EMNLP*, 2014. 6

[17] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 6

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

[19] Rensis Likert. A technique for the measurement of attitudes. *Archives of psychology*, 1932. 8

[20] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. 6

[21] Jingyu Liu, Liang Wang, Ming-Hsuan Yang, et al. Referring expression generation and comprehension via attributes. In *CVPR*, 2017. 7

[22] Li Liu, Wanli Ouyang, Xiaogang Wang, Paul Fieguth, Jie Chen, Xinwang Liu, and Matti Pietikinen. Deep learning for generic object detection: A survey, 2018. 1

[23] Jiasen Lu, Jianwei Yang, Dhruv Batra, and Devi Parikh. Neural baby talk. In *CVPR*, 2018. 2

[24] Ruotian Luo and Gregory Shakhnarovich. Comprehension-guided referring expressions. In *CVPR*, 2017. 2

[25] Minh-Thang Luong, Hieu Pham, and Christopher D Manning. Effective approaches to attention-based neural machine translation. In *EMNLP*, 2015. 2

[26] Junhua Mao, Jonathan Huang, Alexander Toshev, Oana Camburu, Alan L Yuille, and Kevin Murphy. Generation and comprehension of unambiguous object descriptions. In *CVPR*, 2016. 2, 3, 6, 7

[27] Tomáš Mikolov, Martin Karafiát, Lukáš Burget, Jan Černockỳ, and Sanjeev Khudanpur. Recurrent neural network based language model. In *Eleventh Annual Conference of the International Speech Communication Association*, 2010. 1

[28] Varun K Nagaraja, Vlad I Morariu, and Larry S Davis. Modeling context between objects for referring expression understanding. In *ECCV*, 2016. 6, 7

[29] Yulei Niu, Hanwang Zhang, Manli Zhang, Jianhong Zhang, Zhiwu Lu, and Ji-Rong Wen. Recursive visual attention in visual dialog. In *CVPR*, 2019. 2

[30] Jeffrey Pennington, Richard Socher, and Christopher Manning. Glove: Global vectors for word representation. In *EMNLP*, 2014. 6

[31] Bryan A Plummer, Arun Mallya, Christopher M Cervantes, Julia Hockenmaier, and Svetlana Lazebnik. Phrase localization and visual relationship detection with comprehensive image-language cues. In *ICCV*, 2017. 1

[32] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NIPS*, 2016. 1, 2, 6

[33] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. 6

[34] Kai Sheng Tai, Richard Socher, and Christopher D Manning. Improved semantic representations from tree-structured long short-term memory networks. *ACL*, 2015. 4

[35] Andreas Veit and Serge Belongie. Convolutional networks with adaptive inference graphs. In *ECCV*, 2018. 2, 5

[36] Kexin Yi, Jiajun Wu, Chuang Gan, Antonio Torralba, Pushmeet Kohli, and Joshua B Tenenbaum. Neural-symbolic vqa: Disentangling reasoning from vision and language understanding. In *NIPS*, 2018. 2

[37] Licheng Yu, Zhe Lin, Xiaohui Shen, Jimei Yang, Xin Lu, Mohit Bansal, and Tamara L Berg. Mattnet: Modular attention network for referring expression comprehension. In *CVPR*, 2018. 2, 6, 7, 8

[38] Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. Modeling context in referring expressions. In *ECCV*, 2016. 2, 3, 6

[39] Licheng Yu, Hao Tan, Mohit Bansal, and Tamara L Berg. A joint speakerlistener-reinforcer model for referring expressions. In *CVPR*, 2017. 2, 7

[40] Zhou Yu, Jun Yu, Chenchao Xiang, Zhou Zhao, Qi Tian, and Dacheng Tao. Rethinking diversified and discriminative proposal generation for visual grounding. *IJCAI*, 2018. 2

[41] Hanwang Zhang, Zawlin Kyaw, Shih-Fu Chang, and Tat-Seng Chua. Visual translation embedding network for visual relation detection. In *CVPR*, 2017. 1

[42] Hanwang Zhang, Yulei Niu, and Shih-Fu Chang. Grounding referring expressions in images by variational context. In *CVPR*, 2018. 5, 7

# Supplementary Material for "Learning to Assemble Neural Module Tree Networks for Visual Grounding"

## A. Implementation of Tree LSTM

We simplified the implementation of tree LSTM (Eq. 4) in the main paper as:

$$\boldsymbol{c}_t^\uparrow, \boldsymbol{h}_t^\uparrow = \text{TreeLSTM}(\boldsymbol{e}_t, \{\boldsymbol{c}_{tj}^\uparrow\}, \{\boldsymbol{h}_{tj}^\uparrow\}), \quad j \in \mathcal{C}_t, \tag{17}$$

where $\boldsymbol{c}_{tj}^\uparrow$, $\boldsymbol{h}_{tj}^\uparrow$ denote the cell and hidden vectors of the $j$-th child of node $t$. Specifically, our tree LSTM transition equations are:

$$\tilde{\boldsymbol{h}}_t = \sum\nolimits_{j \in \mathcal{C}_t} \boldsymbol{h}_{tj}^\uparrow, \tag{18}$$

$$\boldsymbol{i}_t = \sigma(W^{(i)}\boldsymbol{e}_t + U^{(i)}\tilde{\boldsymbol{h}}_t + \boldsymbol{b}^{(i)}), \tag{19}$$

$$\boldsymbol{o}_t = \sigma(W^{(o)}\boldsymbol{e}_t + U^{(o)}\tilde{\boldsymbol{h}}_t + \boldsymbol{b}^{(o)}), \tag{20}$$

$$\boldsymbol{u}_t = \tanh(W^{(u)}\boldsymbol{e}_t + U^{(u)}\tilde{\boldsymbol{h}}_t + \boldsymbol{b}^{(u)}), \tag{21}$$

$$\boldsymbol{f}_{tj} = \sigma(W^{(f)}\boldsymbol{e}_t + U^{(f)}\boldsymbol{h}_{tj}^\uparrow + \boldsymbol{b}^{(f)}), \tag{22}$$

$$\boldsymbol{c}_t^\uparrow = \boldsymbol{i}_t \odot \boldsymbol{u}_t + \sum\nolimits_{j \in \mathcal{C}_t} \boldsymbol{f}_{tj} \odot \boldsymbol{c}_{tj}^\uparrow, \tag{23}$$

$$\boldsymbol{h}_t^\uparrow = \boldsymbol{o}_t \odot \tanh(\boldsymbol{c}_t^\uparrow), \tag{24}$$

where $\odot$ is the element-wise multiplication, $\sigma(\cdot)$ is the sigmoid function, $W$, $U$, $b$ are trainable parameters.

## B. More Qualitative Results

In this section, we provide more qualitative results to demonstrate the internal reasoning steps of NMTREE. In Figure 8, we visualize the reasoning process inside `Comp` modules. In Figure 9, Figure 10, and Figure 11, we visualize the tree structures, the module assembly, the attention map at each intermediate step, and the final results. Specifically, Figure 9 are qualitative results with ground-truth bounding boxes. As comparison, we also show some failure cases. Figure 10 are qualitative results with detected bounding boxes. Figure 11 are qualitative results with detected masks.
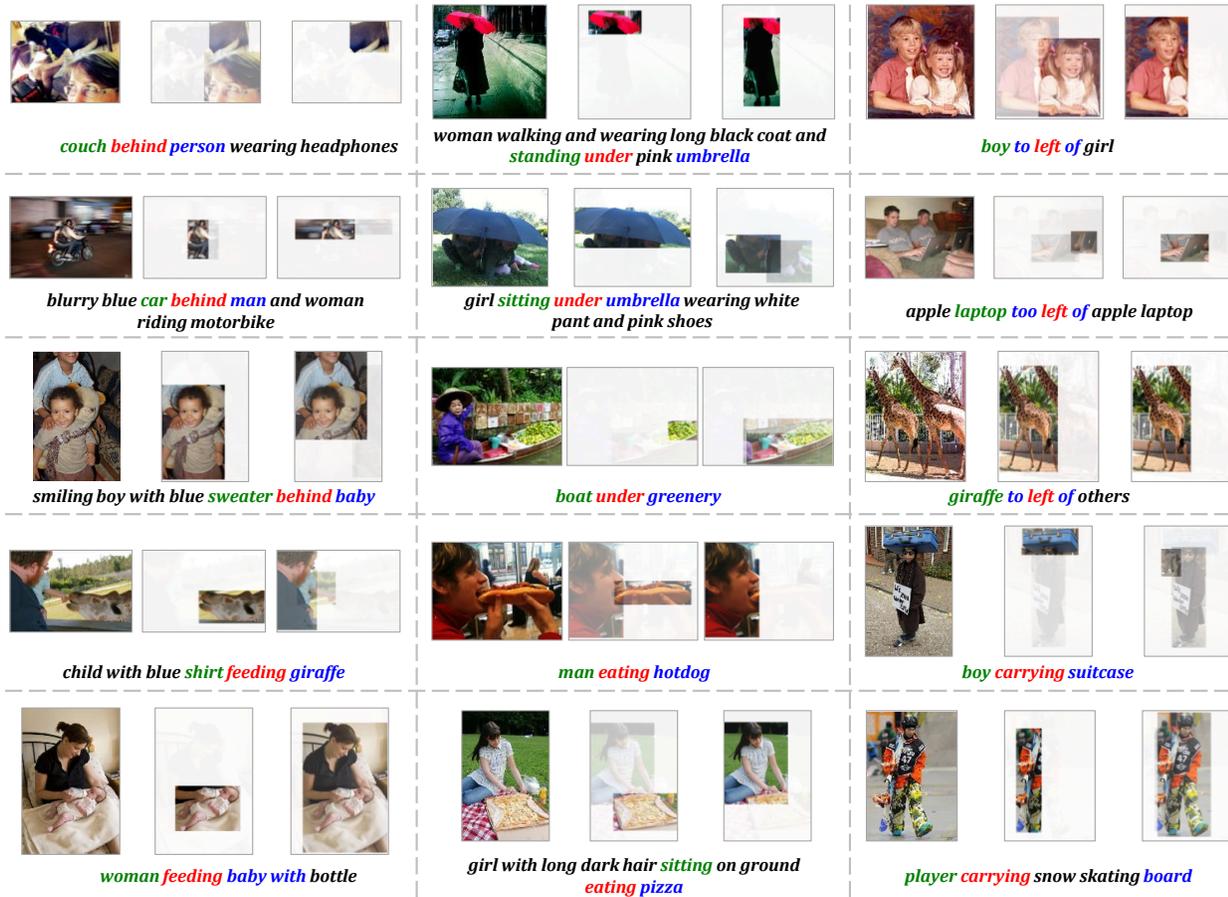
Figure 8. The compositional reasoning inside Comp. Each example contains the original image (left), the contextual attention map (middle) and the output attention map (right). We represent partial tree structure by colors: red for the current node, blue for children and green for parent.

(a) a chair behind a woman drinking wine

(b) woman standing looking the rails at the ocean with a man seated on a brown bench

(c) the empty chair in front of the laptop

(d) a blue car behind the man walking an elephant

(a) correct

(e) a brown chair with a guy sitting on it

(f) the left ski of a skier

(g) a chair that is behind the laptop

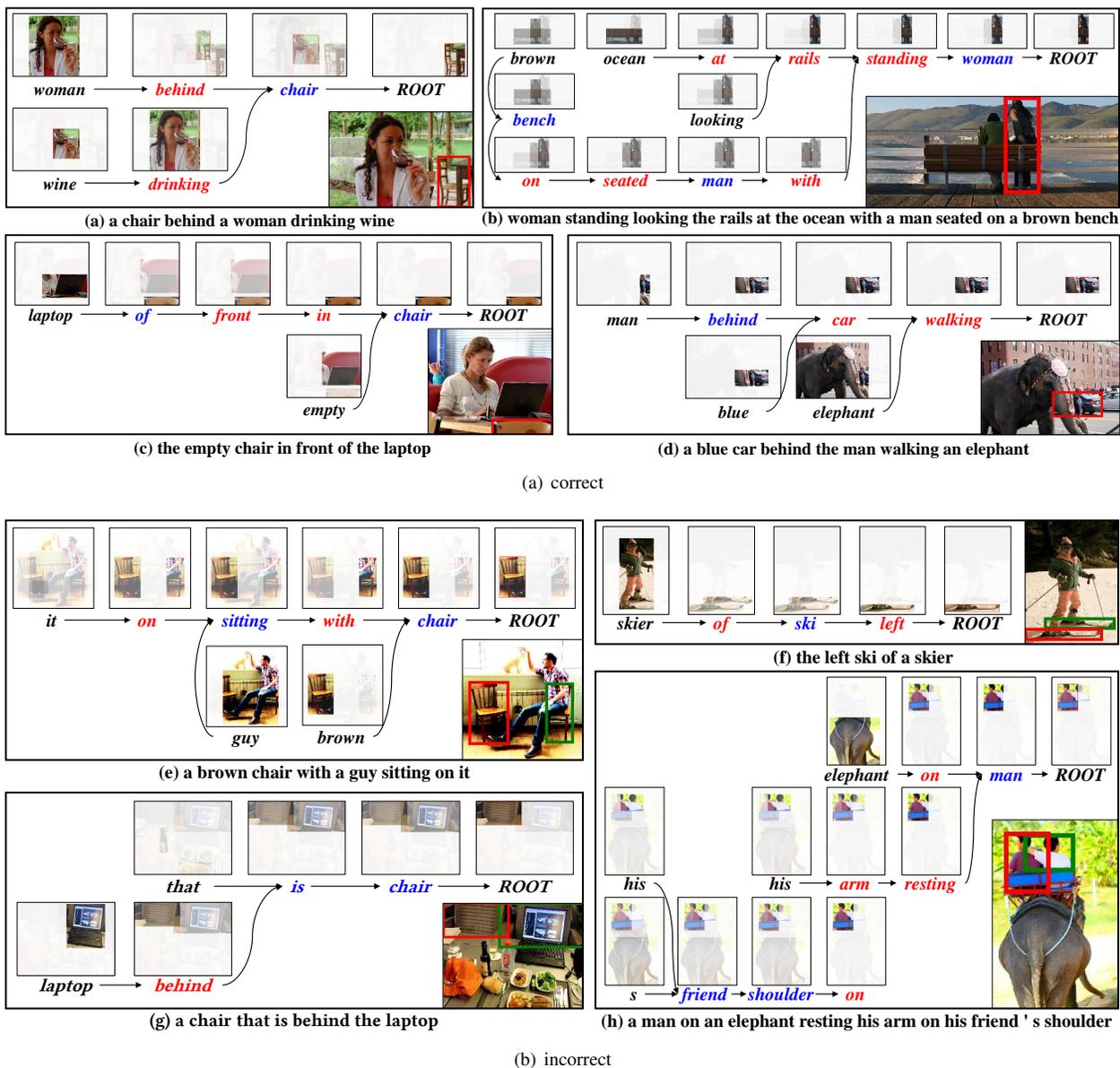(h) a man on an elephant resting his arm on his friend ' s shoulder

(b) incorrect

Figure 9. Qualitative results with ground-truth bounding boxes. Words in different colors indicate corresponding modules: black for Single, red for Comp, and blue for Sum. The bottom right corner is the original image with a green bounding box as ground-truth and a red bounding box as our result.

**(a) a man with an orange t - shirt with pizza**

**(b) an adult elephant standing with a baby elephant**

**(c) a blue color car parked near a tree**

**(d) man leaning on bike on boat**

Figure 10. Qualitative results with detected bounding boxes. Words in different colors indicate corresponding modules: black for `Single`, red for `Comp`, and blue for `Sum`. The bottom right corner is the original image with a green bounding box as ground-truth and a red bounding box as our result.
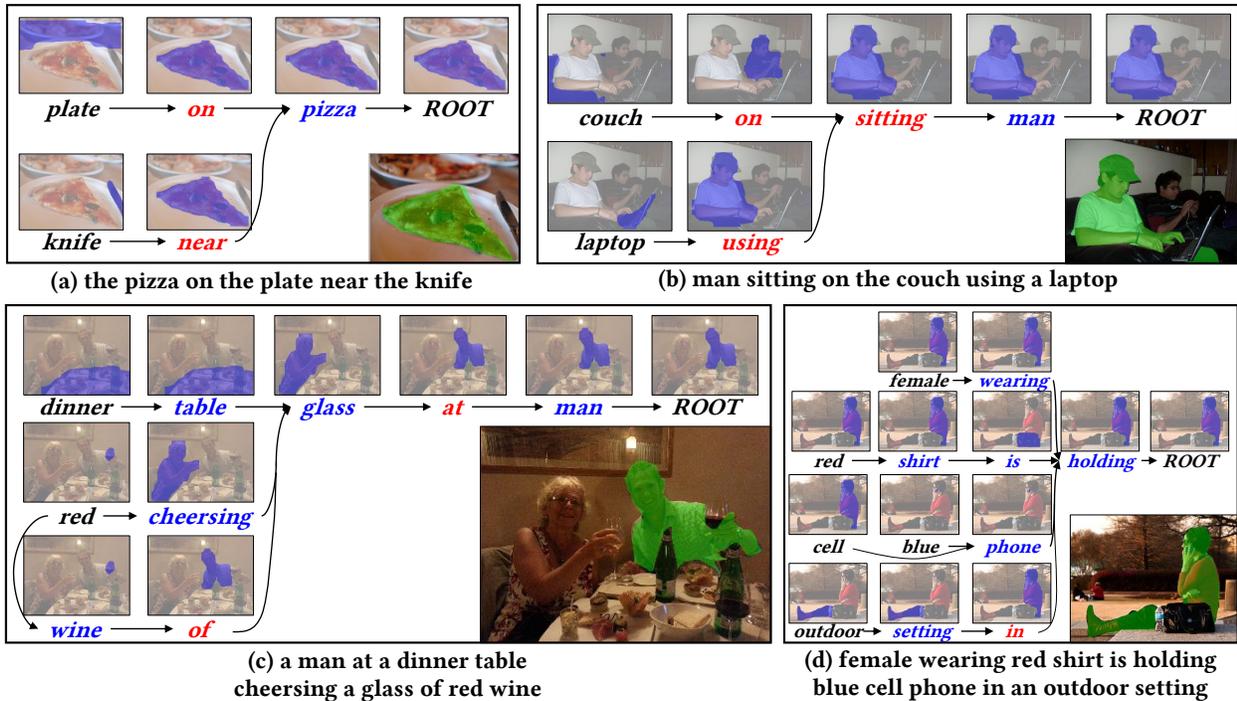


**(a) the pizza on the plate near the knife**

**(b) man sitting on the couch using a laptop**

**(c) a man at a dinner table cheersing a glass of red wine**

**(d) female wearing red shirt is holding blue cell phone in an outdoor setting**

Figure 11. Qualitative results with detected masks. Words in different colors indicate corresponding modules: black for `Single`, red for `Comp`, and blue for `Sum`. The blue masks indicate the regions with maximum score, and the green masks indicate the ground-truth.