

Supplementary Material for “Learning to Assemble Neural Module Tree Networks for Visual Grounding”

Daqing Liu¹ Hanwang Zhang² Feng Wu¹ Zheng-Jun Zha¹

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

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

1. Implementation of Tree LSTM

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

$$\mathbf{c}_t^\uparrow, \mathbf{h}_t^\uparrow = \text{TreeLSTM}(\mathbf{e}_t, \{\mathbf{c}_{tj}^\uparrow\}, \{\mathbf{h}_{tj}^\uparrow\}), \quad j \in \mathcal{C}_t, \quad (1)$$

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

$$\tilde{\mathbf{h}}_t = \sum_{j \in \mathcal{C}_t} \mathbf{h}_{tj}^\uparrow, \quad (2)$$

$$\mathbf{i}_t = \sigma(W^{(i)}\mathbf{e}_t + U^{(i)}\tilde{\mathbf{h}}_t + \mathbf{b}^{(i)}), \quad (3)$$

$$\mathbf{o}_t = \sigma(W^{(o)}\mathbf{e}_t + U^{(o)}\tilde{\mathbf{h}}_t + \mathbf{b}^{(o)}), \quad (4)$$

$$\mathbf{u}_t = \tanh(W^{(u)}\mathbf{e}_t + U^{(u)}\tilde{\mathbf{h}}_t + \mathbf{b}^{(u)}), \quad (5)$$

$$\mathbf{f}_{tj} = \sigma(W^{(f)}\mathbf{e}_t + U^{(f)}\mathbf{h}_{tj}^\uparrow + \mathbf{b}^{(f)}), \quad (6)$$

$$\mathbf{c}_t^\uparrow = \mathbf{i}_t \odot \mathbf{u}_t + \sum_{j \in \mathcal{C}_t} \mathbf{f}_{tj} \odot \mathbf{c}_{tj}^\uparrow, \quad (7)$$

$$\mathbf{h}_t^\uparrow = \mathbf{o}_t \odot \tanh(\mathbf{c}_t^\uparrow), \quad (8)$$

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

2. More Qualitative Results

In this section, we provide more qualitative results to demonstrate the internal reasoning steps of NMTree. In Figure 1, we visualize the reasoning process inside Comp modules. In Figure 2, Figure 3, and Figure 4, we visualize the tree structures, the module assembly, the attention map at each intermediate step, and the final results. Specifically, Figure 2 are qualitative results with ground-truth bounding boxes. As comparison, we also show some failure cases. Figure 3 are qualitative results with detected bounding boxes. Figure 4 are qualitative results with detected masks.

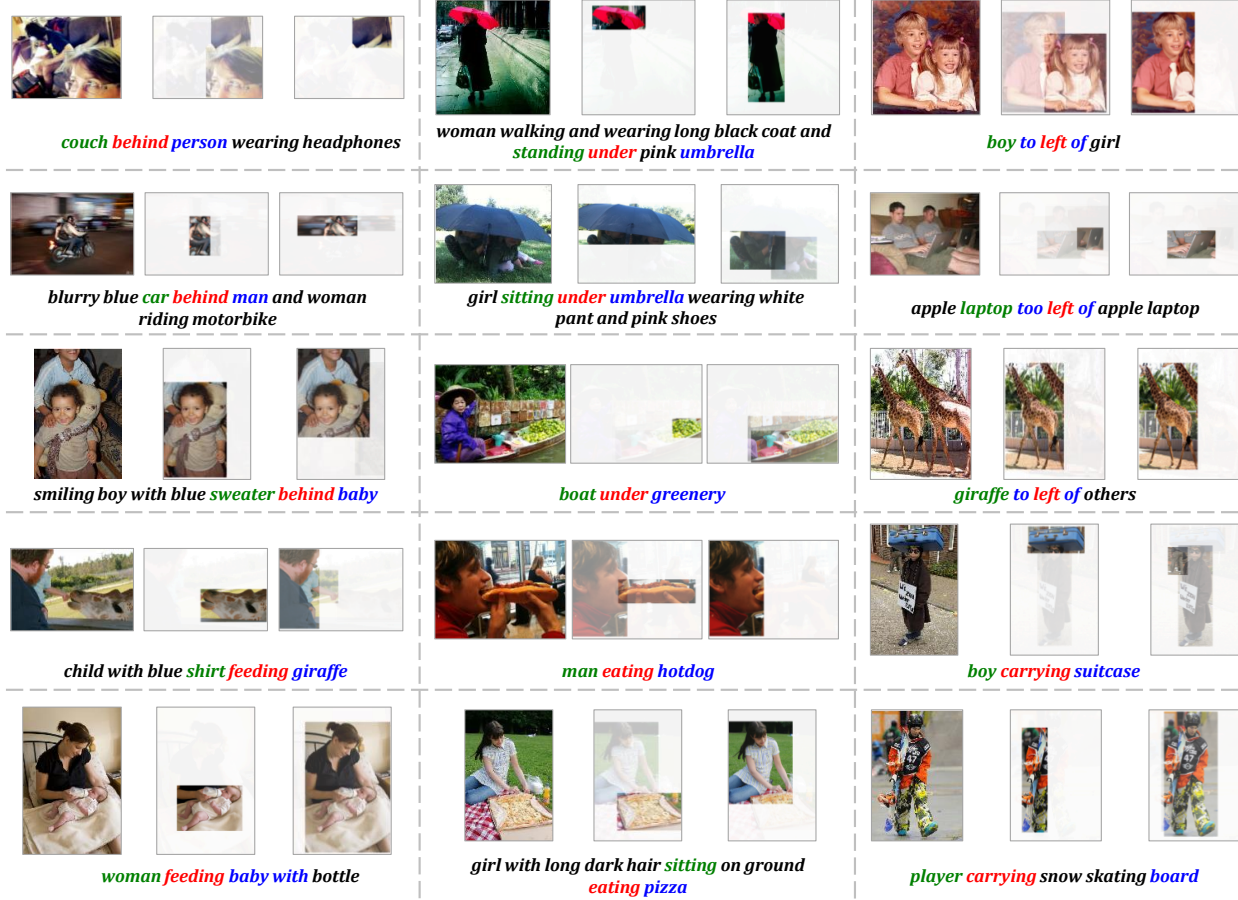
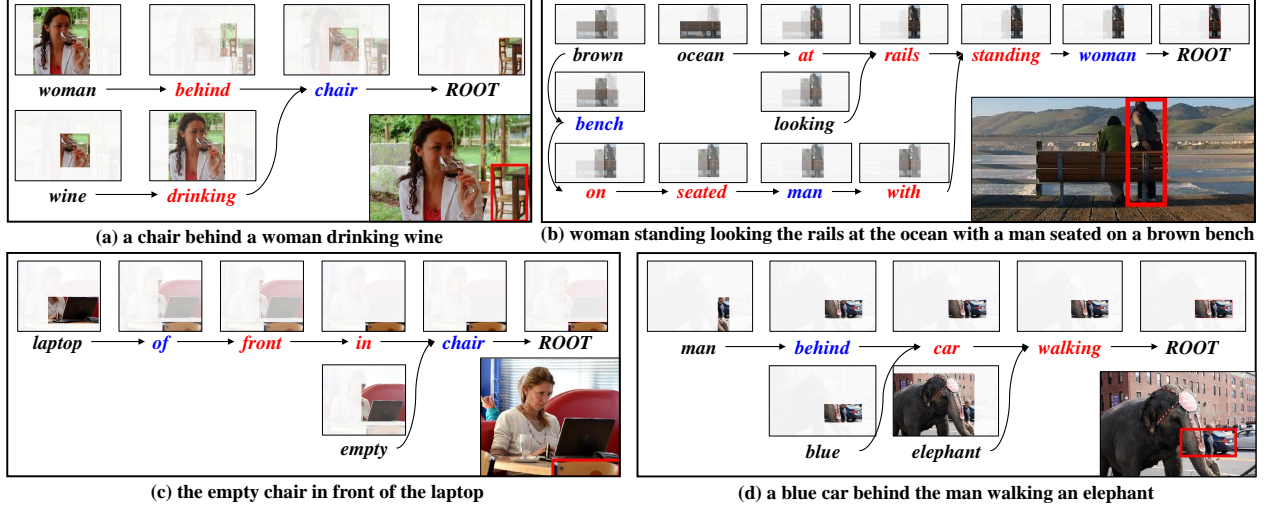
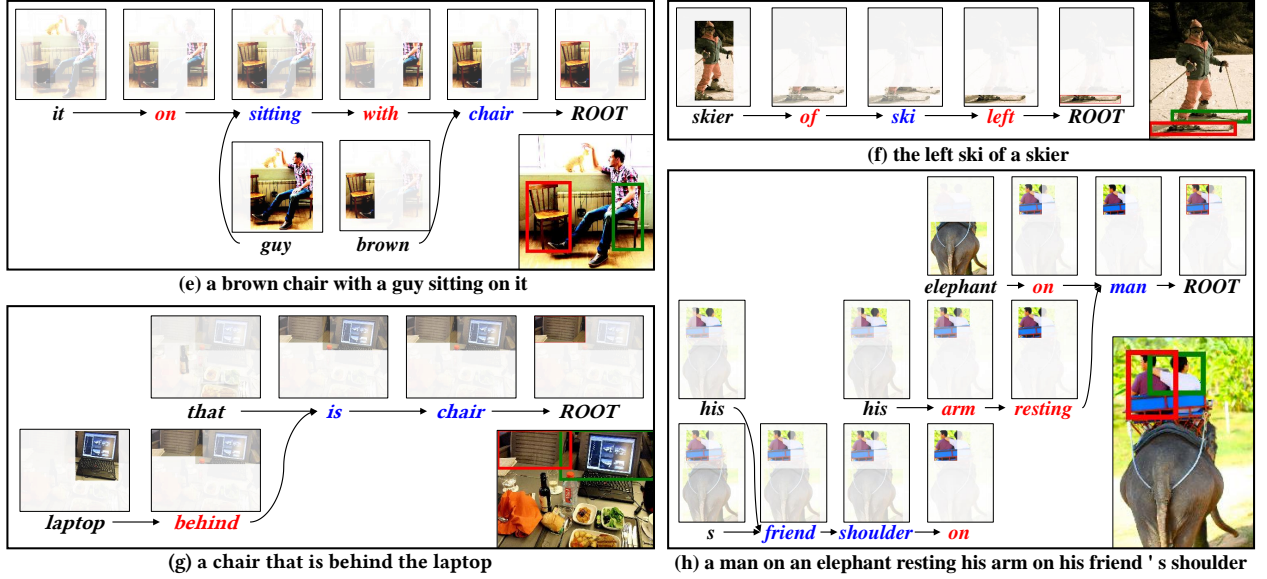


Figure 1. The compositional reasoning inside Comp. Each example contains the original image (left), the contextual attention map (middle) and the output attention map (right). We represent partial tree structure by colors: red for the current node, blue for children and green for parent.



(a) correct



(b) incorrect

Figure 2. 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.

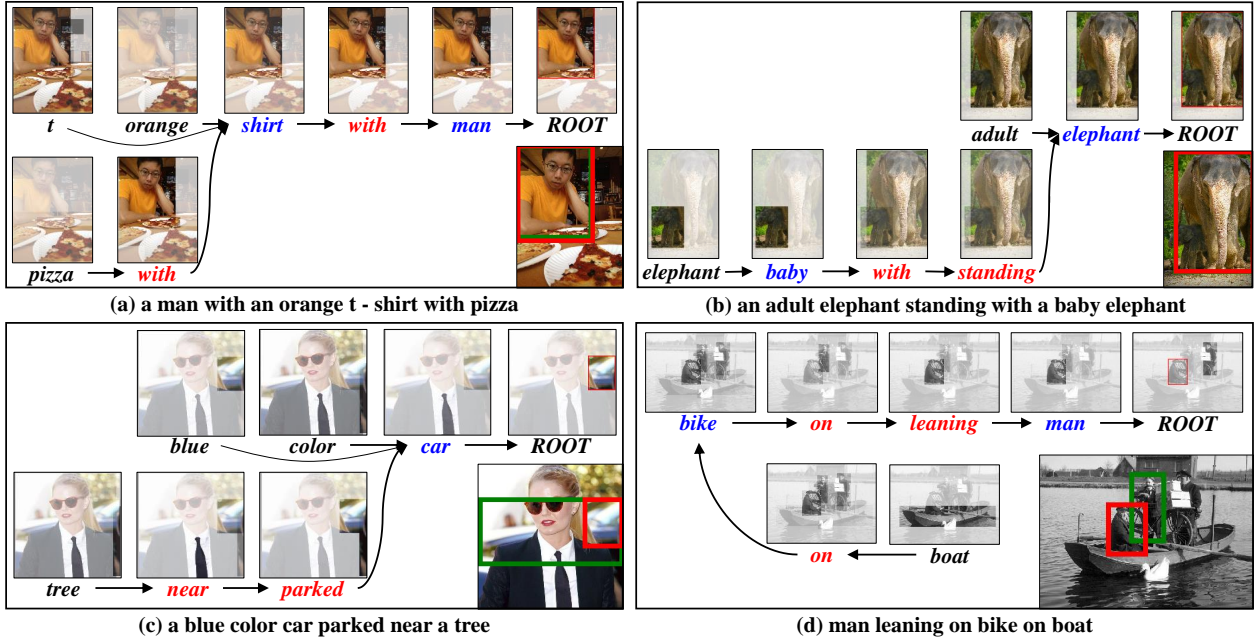


Figure 3. Qualitative results with detected bounding boxes. Words in different colors indicate corresponding modules: black for Single, red for Comp, and blue for Sum. The bottom right corner is the original image with a green bounding box as ground-truth and a red bounding box as our result.

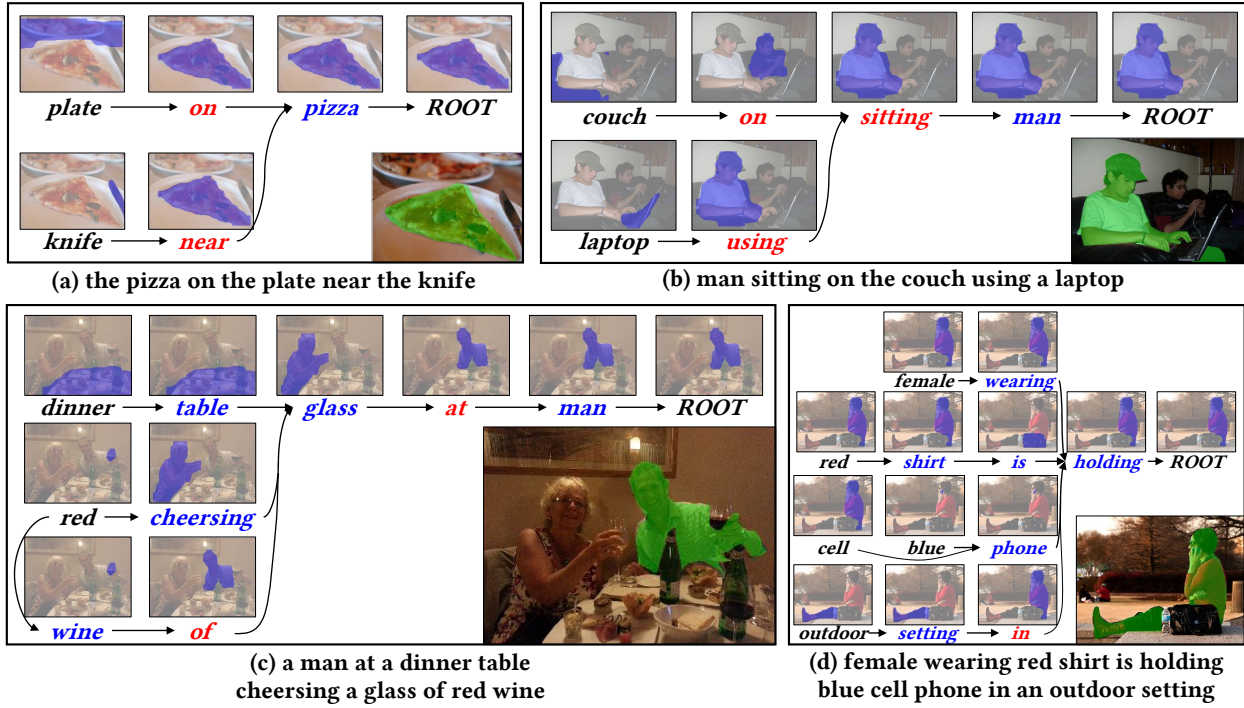


Figure 4. Qualitative results with detected masks. Words in different colors indicate corresponding modules: black for Single, red for Comp, and blue for Sum. The blue masks indicate the regions with maximum score, and the green masks indicate the ground-truth.