# Supplementary Material for "Improving One-stage Visual Grounding by Recursive Sub-query Construction"

Zhengyuan Yang[1]    Tianlang Chen[1]    Liwei Wang[2]    Jiebo Luo[1]

[1]University of Rochester    [2]Tencent AI Lab, Bellevue
{zyang39,tchen45,jluo}@cs.rochester.edu, liweiwang@tencent.com

## 1  Qualitative Results

We present additional qualitative results in Figures A, B to visualize the recursive disambiguation process. We highlight the following scenarios.

- **Recursive disambiguation.** The proposed recursive sub-query construction framework improves one-stage visual grounding by addressing the current limitations on grounding long queries. We observe a desired recursive disambiguation process that the text-conditional visual feature step by step generates more accurate and confident predictions.
  As an easy case, better modeling the modifiers of the head noun already corrects a portion of previous failures. For example, in Figure A (a), the peak in the heatmap moves from the tennis player to the referred person in the back, after observing the modifier "watching" in the second round. On the contrary, previous one-stage methods tend to overemphasize the head nouns, without full consideration of the modifiers. Our proposed method also works well on more complex queries via recursive disambiguation, such as in the example "persons head with drill in the middle" from the main paper, and "pony tail lady on the right forefront" in Figure A (b).
- **Challenging regions.** We observe that our method performs well on challenging examples such as Figures A (c) and (d), where the referred target is tiny, the scene contains visually similar distracting objects, and the query includes complex attributes and relationship descriptions.
- **Attributes.** Figures A (e) and (f) include examples that require the correct understanding of attributes such as color and size. For example, in the final round of Figure A (e), the peak moves from the distracting object "trolley" to the referred red suitcase by observing the constructed sub-query "red."
- **Failure cases.** Figures A (g) and (h) show failure cases of our model. The model either misses certain related objects such as the "controller" in Figures A (g) or fails to understand some rarely appeared attributes such as "plaid" in Figures A (h). Therefore, one or more distracting objects remain to have high heatmap responses in the final round, despite the model might still predict the correct bounding box.

We present additional qualitative results in Figure B.

**Table A.** Ablation studies on number of rounds.

| #Rounds | Acc@0.5 | Time(ms) |
|---|---|---|
| $K = 1$ | 58.22 | 23 |
| $K = 2$ | 59.31 | 25 |
| $K = 3$ (Ours) | 60.96 | 26 |
| $K = 4$ | 61.08 | 28 |
| $K = 5$ | 60.80 | 30 |
| $K = 6$ | 61.00 | 32 |

**Table B.** Ablation studies on the modifications in Ours-Large.

| Method | Acc@0.5 | Time(ms) |
|---|---|---|
| Ours-Base | 60.96 | 26 |
| + ConvLSTM | 61.26 | 27 |
| + Size 512 | 61.99 | 34 |
| + Both | 63.12 | 36 |

**Table C.** Performance break-down with attributes.

| *ReferItGame* | Color | Loc. | Size | All |
|---|---|---|---|---|
| Percent | 7.84 | 53.63 | 7.00 | 100.00 |
| One-Stage-BERT | 43.07 | 50.98 | 53.83 | 59.30 |
| Ours-Base | 50.52 | 58.71 | 61.83 | 64.33 |
| **Relative Gain** | 17.30 | 15.16 | 14.86 | 8.48 |

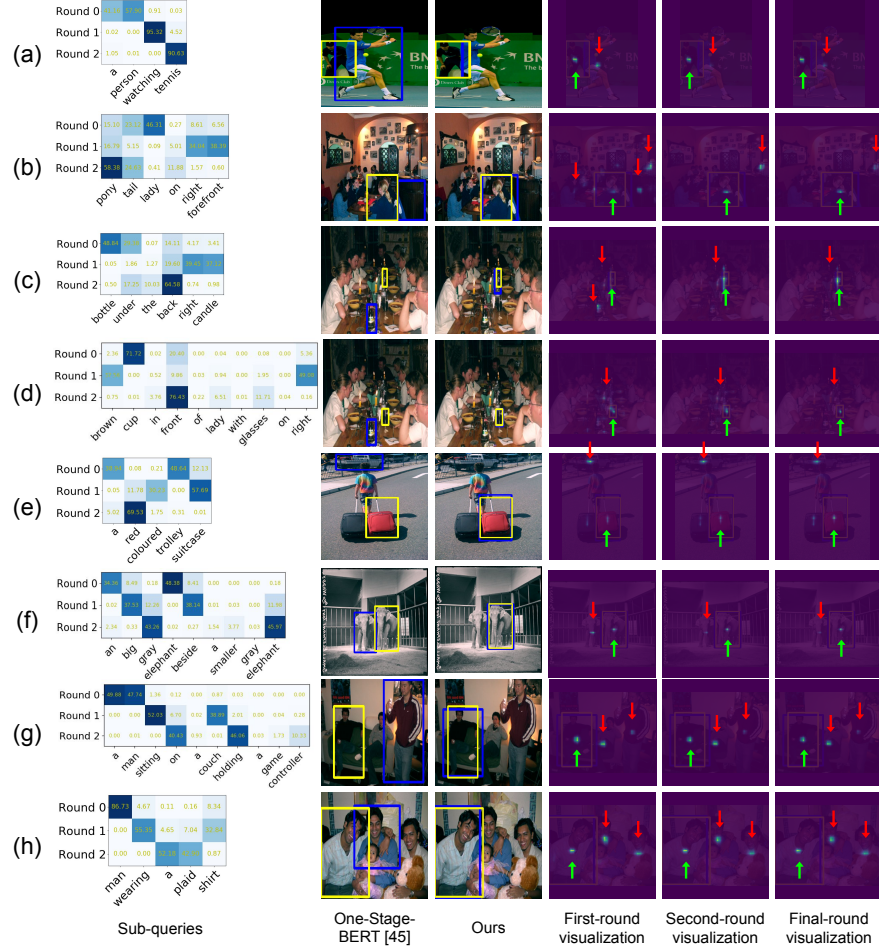| *RefCOCOg* | Color | Loc. | Size | All |
|---|---|---|---|---|
| Percent | 18.54 | 32.10 | 12.48 | 100.00 |
| One-Stage-BERT | 55.06 | 56.46 | 57.85 | 58.70 |
| Ours-Base | 62.92 | 66.90 | 66.53 | 64.87 |
| **Relative Gain** | 14.28 | 18.49 | 15.00 | 10.51 |

## 2   Ablation Studies

**Number of rounds.** Table A shows the ablation studies on the different number of rounds $K$ on the RefCOCOg-google dataset. We observe that increasing the number of rounds does not lead to an increase in accuracy after a dataset-specific threshold (*e.g.*, $K \geq 3$ on RefCOCOg). Therefore, we select $K = 3$ as the default value in our experiments for a balance between efficiency and accuracy. Although we report all results with $K = 3$ in the main paper's Table 1 for clarity, we note that a different $K$ might slightly improve the accuracy or reduce the inference time on different datasets.

**Ours-Large.** We observe several modules and settings that further improve the grounding accuracy, but meanwhile slightly slow the inference speed or increase the model complexity. Therefore, we list such modifications separately and refer to the corresponding framework "Ours-Large." As shown in Table B, we observe an increase in accuracy with a larger input image size. Furthermore, we improve the accuracy by using all intermediate text-conditional visual feature with a ConvLSTM module. The ConvLSTM module takes feature $\{v^{(k)}\}_{k=1}^{K}$ as the input, and outputs the last hidden state for grounding box prediction.

**Performance break-down with attributes.** The performance break-down study in Section 4.4 shows the effectiveness of our proposed method in modeling and grounding long queries. Other than the improvements on long queries, our method also shows advantages in modeling queries with attribute descriptions, such as color, location, or size. To validate the observation, we construct the attribute subsets from the test set of ReferItGame and RefCOCOg based on the contained attribute keywords, *e.g.*, *"white," "black," "red," "blue,"* etc., for "color;" *"right," "left," "front," "middle,"* etc., for "location;" *"big," "little," "small," "tall,"* etc., for "size."

The first row of Table C shows the experimented dataset and the name of the subset. "Color," "location," and "size" indicate that the query in the subset contains at least one corresponding attribute keywords. "All" reports the

performance on the entire dataset. The second row shows the portion of samples in each subset. The remaining rows indicate the grounding accuracy and the relative gain. As shown in the last row, the relative gains on the attribute subsets are around 15% and are higher than the relative gain on the entire dataset of around 10% (*cf.* the middle three and the last column of Table C).



**Fig. A.** Visualizations of the constructed sub-queries and the intermediate text-conditional visual feature at each round. Blue/ yellow boxes are the predicted regions/ ground truths. The green up arrow and the red down arrow highlight the target and the major distracting objects on heatmaps, respectively.

**Fig. B.** Additional qualitative results.