1 We would like to thank all reviewers for your constructive and intriguing comments. We have concluded more related
2 work, compared with more SOTA, based on your suggestions. We have significantly polished our paper in terms of
3 introduction reorganization, grammar correction, fluent expression as well as broader impacts for camera ready.

**Reviewer #1 and Reviewer #4**

5 Q1: Comparison with FiLM (AAAI 2018). The experiment is somewhat inadequate.

6 A1: Feature-wise Linear Modulation, namely FiLM is proposed to influence neural network computation via a simple,
7 feature-wise affine transformation based on conditioning information. Such a novel module has achieved SOTA
8 performance on the CLEVR benchmark. The GRU module in FiLM functions similar to the BiLSTM module in our
9 porposed SIRI, and the FiLM functions similar to the spatial relation guided distillation module. Besides, positional
10 embedding are leveraged in both work. Whereas, the spatial relation guided distillation we introduced induces the
11 network with a gate mechanism while FiLM performs feature-wise affine transformation. To fairly compare our SIRI
12 with FiLM, we adopt the part II and the part III only in SIRI. Our SIRI has an accuracy@80px of 58.33%, much better
13 than FiLM with an accuracy@80px of 52.37%, which demonstrates the effectiveness of our SIRI.

**Reivewer #2**

15 Q1: The effect of k (for the number of top-k selected orientation words) is also not studied.

16 A1: The accuracies@80px are 51.02%, 58.33% and 59.74% when k is 4, 6 and 8, respectively.

17 Q3: The ablation study is not as thorough as it can be (it adds the components in order). Ideally, it would also show the
18 effect of 2) and 3) without using 1) (with ResNet features instead of GloRe) and effect of 3) without 2).

19 A3: The accuracy@80px is 56.24% when the component 2) and 3) are adopted without 1). Besides, the accuracy is
20 38.76% when the component 3) is adopted without 1) and 2).

21 Q4: It is not clear whether it would be of interest to the broader NeurIPS community.

22 A4: We believe that our proposed SIRI with strong novelty and promising performance provides a new insight in terms
23 of architecture design for any vision-language tasks, such as VQA, SDR, etc.

**Reviewer #3**

25 Q2: The new extended dataset used in 4.2 should be described more, no details are given on this new data and how it
26 was collected.

27 A2: Due to the page limitation, the details of the new extended dataset are appended in the supplementray materials.
28 We also analyze the word frequency on it and visualize the prediction results on this new dataset.

**Reviewer #4**

30 Q2: Although this is a new task, but the solution is compositional and of limited novelty.

31 A2: In this paper, we design a novel framework to explicitly tackle the SDR task. Each component is carefully designed
32 and well investigated. We believe that such a novel framework can push forward this important task and provides a new
33 insight for any other vision-language task.

34 Q3: The ablation study in Table 3 only shows three combinations besides pure LingUnet – I, I+II and I+II+III, it will be
35 better if the author could provide the combinations of II + III, III only and II only. Then we can better evaluate the
36 importance of Part II and Part III.

37 A3: The accuracies@80px are 56.24%, 38.76% and 44.95% for II + III, III only, II only, respectively.

38 Q4: About the generalization ability of Stage II & III.

39 A4: Because MAttNet detects objects and scores the RoIs, the Stage II and III in our paper are not suitable to MAttNet.
40 We use YOLO-VG(A Fast and Accurate One-Stage Approach to Visual Grounding, ICCV2019) as our baseline, a one
41 stage method for visual grounding, which can be end-to-end trained. To investigate the generalization ability of Stage II
42 & III, we add the Stage II and III to it. The results show that Stage II and III can improve the performance by 0.9% and
43 0.7%, respectively. Thus, our proposed modules can also perform well on other tasks and datasets.

44 Q5: The paper has some typos, such as: [Line 92]: adapt should be adopt; [Line 121]: averaged should be summed; [In
45 figure 2]: VI (in fact 6) should be IV (4).

46 A5: Thanks for pointing out these typos. We have significantly polished this paper for camera ready.