| Model | $Q \rightarrow A$ | | $QA \rightarrow R$ | | $Q \rightarrow AR$ | |
|---|---|---|---|---|---|---|
| | GloVe | BERT | GloVe | BERT | GloVe | BERT |
| BottomUpTopDown [2] | 42.8 | 62.3 | 25.1 | 63.0 | 10.7 | 39.6 |
| R2C [46] | 46.4 | 63.8 | 38.3 | 67.2 | 18.3 | 43.1 |
| **HGL (ours)** | **54.1** | **69.4** | **42.7** | **70.6** | **25.1** | **49.1** |

Table 1: Results on the VCR validation set. While BERT helps the performance of these baselines, our model still performs the best.

1 **To All:** We thank all reviewers for the constructive and insightful comments and address the comments as follows. We
2 appreciate the fact that all reviewers give positive evaluation results about our HGL. We promise to revise typos and
3 grammatical mistakes and polish our paper in the revised version based on these valuable suggestions.

4 **Q1: Clarification on results of Fig.4 and Fig.5 (R1).** In Fig.4(a), "feeling" from question can be aligned with the
5 most suitable semantic word "anger" from the answer for right answer prediction, which demonstrated the effectiveness
6 of our QAHG. In Fig.4(e), the visual representation "person5" is aligned with the linguistic word "witness" correctly
7 via VAHG, because "person5" is the "witness" in this scenario. Visual representation "person1" (tag associated with the
8 red box is fed to model) is connected with the emotional word "angry" to achieve a heterogeneous relationship. In Fig.5,
9 the person would fell wet because of rain in the scenario. The arrow is pointing at raindrop.

10 **Q2: Comment on CVM applied to vision (R1).** The CVM is suitable to apply to visual context, because there
11 are more information can be obtained from the visual evidence. For instance, the first example of Figure 1 from
12 supplementary material, the motivation of person1 must get information from visual evidence (e.g. cake1, table and
13 background) instead of the question to predict the right answer and reason.

14 **Q3: Discussions on the affection of concatenation in Eq.5 (R1).** To combine the relationship between $X_m$ and $Y_o$
15 with $Y_o$ instead of simply combining $Y_o$ with $X_m$, we need the concatenation operation in Eq.5. We did experiments to
16 analyze the affection and found the concatenation can improve ~0.4% accuracy (69 vs 69.4) on Q->A task over VCR
17 validation set. The experiment showed that the concatenation operation can help the model to further understand the
18 heterogeneous relationship for predicting correctly.

19 **Q4: More experiments on common VQA datasets (R2).** We trained our HGL on VQA2 dataset following the same
20 experimental setup [2], and evaluated our HGL on VQA2 validation dataset compared with common VQA approach
21 such as BottomUpTopDown [2]. The BottomUpTopDown [2] achieved 63.2 accuracy across all question types and the
22 HGL achieved 65.3 accuracy, which showed our proposed model can also be transferred to VQA task and perform well.

23 **Q5: Ablation experiments between BERT and common VQA backbone (e.g. GloVe) (R2).** In Table 1, the BERT
24 indeed helped the performance of HGL ($54.1 \rightarrow 69.4$, $42.7 \rightarrow 70.6$, $25.1 \rightarrow 49.1$), our HGL still got SOTA performance
25 compared with common VQA approach such as BottomUpTopDown [2], and R2C [46] on VCR validation set.

26 **Q6: Clarification on CVM and the possibility to take "visual context" as graph nodes (R2).** The CVM aims to
27 capture ambiguous semantic context (e.g. rainy/snowy weather) that lack of specific labels for detection and can
28 not benefit from the labeled object grounding boxes and categories such as "person" and "dog" during training. The
29 long-range visual context means spacial visual evidence that is consisted of pixel-wise representations of different
30 positions. It may not be suitable to take "visual context" as graph nodes. Because the "visual context" contains complex
31 and uncertain semantics such as different weather and emotional expression, which is different to make sure each node
32 semantics and the number of nodes to represent the whole visual context.

33 **Q7: Clarification on graph learning, inputs and outpus of the model during training and inference (R3).** The
34 graphs are latent to learn. The inputs of our HGL are candidate answers, a question and an image. The output of the
35 model is a prediction of right answers. The inputs and outputs of the model are the same during training and inference.

36 **Q8: Explanation for a few questions on section 1 and section 3 (R1 R2 R3).** In section 1, the concept of information
37 isolated island in our paper refers to the independent of different semantic nodes can not achieve semantic inference in
38 a homogeneous graph that connects similar semantic nodes by attribute (e.g. Figure 1(a)) or grammar (e.g. Figure 1(b)).
39 The $y_i^{rl}$ should be $y_i^l$ in Eq.12, we will revise the typos in final version. In L131, the correlations between the vectors are
40 computed via a dot product, we promise to give a more precise definition of this operation in the final version. In section
41 3.1, maybe "Graph Construction" would be better than the title "Task Definition", we also realize the inappropriate title
42 and will consider carefully the name of the title in our final version. We will briefly describe the three tasks (Q->A,
43 QA->R and Q->AR) and some specific terminologies in our final version to make the paper easier to understand.

44 **Q9: More qualitative analysis and the application of this work (R3).** Due to the space limit, we will show more
45 qualitative results with analysis in our final version. Our HGL can seamlessly integrate the intra-graph and inter-graph
46 reasoning in order to bridge vision and language domain and can interactively refine reasoning paths for semantic
47 agreement. Such capability would show more potentials in facilitating high-level applications that require cross-domain
48 semantic alignment among visual concepts and linguistic words. (e.g. visual grounding and VQA).