# REXUP: I REason, I EXtract, I UPdate with Structured Compositional Reasoning for Visual Question Answering

Siwen Luo*, Soyeon Caren Han*✉, Kaiyuan Sun, and Josiah Poon

School of Computer Science, The University of Sydney, 1 Cleveland Street,
The University of Sydney, NSW 2006, Australia,
{siwen.luo, caren.han, kaiyuan.sun, josiah.poon}@sydney.edu.au

**Abstract.** Visual Question Answering (VQA) is a challenging multi-modal task that requires not only the semantic understanding of images and questions, but also the sound perception of a step-by-step reasoning process that would lead to the correct answer. So far, most successful attempts in VQA have been focused on only one aspect; either the interaction of visual pixel features of images and word features of questions, or the reasoning process of answering the question of an image with simple objects. In this paper, we propose a deep reasoning VQA model (REXUP- REason, EXtract, and UPdate) with explicit visual structure-aware textual information, and it works well in capturing step-by-step reasoning process and detecting complex object-relationships in photo-realistic images. REXUP consists of two branches, image object-oriented and scene graph-oriented, which jointly works with the super-diagonal fusion compositional attention networks. We evaluate REXUP on the benchmark GQA dataset and conduct extensive ablation studies to explore the reasons behind REXUPs effectiveness. Our best model significantly outperforms the previous state-of-the-art, which delivers 92.7% on the validation set, and 73.1% on the test-dev set.

**Keywords:** GQA · Scene Graph · Visual Question Answering

## 1 Introduction

Vision-and-language reasoning requires the understanding and integration of visual contents and language semantics and cross-modal alignments. Visual Question Answering (VQA) [2] is a popular vision-and-language reasoning task, which requires the model to predict correct answers to given natural language questions based on their corresponding images. Substantial past works proposed VQA models that focused on analysing objects in photo-realistic images but worked well only for simple object detection and yes/no questions [17,25,14]. To overcome this simple nature and improve the reasoning abilities of VQA models, the Clever dataset[13] was introduced with compositional questions

---

* Both authors are first author

and synthetic images, and several models [9,20] were proposed and focused on models' inferential abilities. The state-of-the-art model on the Clevr dataset is the compositional attention network(CAN)[11], which generates reasoning steps attending over both images and language-based question words. However, the Clevr dataset is specifically designed to evaluate reasoning capabilities of a VQA model. Objects in the Clevr dataset images are only in three different shapes and four different spatial relationships, which results in simple image patterns. Therefore, a high accuracy on Clevr dataset hardly prove a high object detection and analysis abilities in photo-realistic images, nor the distinguished reasoning abilities of a VQA model. To overcome the limitations of VQA and Clevr [2,7], the GQA dataset [12] includes photo-realistic images with over 1.7K different kinds of objects and 300 relationships. GQA provides diverse types of answers for open-ended questions to prevent models from memorizing answer patterns and examine the understanding of both images and questions for answer prediction.

The state-of-the-art models on the Clevr and VQA dataset suffered large performance reductions when evaluated on GQA [11,19,1]. Most VQA works focus on the interaction between visual pixel features extracted from images and question features while such interaction does not reflect the underlying structural relationships between objects in images. Hence, the complex relationships between objects in real images are hard to learn.Inspired by this motivation, we proposed REXUP(REason, EXtract, UPdate) network to capture step-by-step reasoning process and detect the complex object-relationships in photo-realistic images with the scene graph features. A scene graph is a graph representation of objects, attributes of objects and relationships between objects where objects that have relations are connected via edge in the graph.

The REXUP network consists of two parallel branches where the image object features and scene graph features are respectively guided by questions in an iterative manner, constructing a sequence of reasoning steps with REXUP cells for answer prediction. A super-diagonal fusion is also introduced for a stronger interaction between object features and question embeddings. The branch that processes scene graph features captures the underlying structural relationship of objects, and will be integrated with the features processed in another branch for final answer prediction. Our model is evaluated on the GQA dataset and we used the official GQA scene graph annotations during training. To encode the scene graph features, we extracted the textual information from the scene graph and used Glove embeddings to encode the extracted textual words in order to capture the semantic information contained in the scene graph. In the experiments, our REXUP network achieved the state-of-the-art performance on the GQA dataset with complex photo realistic images in deep reasoning question answering task.

## 2    Related Work and Contribution

We explore research trends in diverse visual question answering models, including fusion-based, computational attention-based, and graph-based VQA models.

**Fusion-based VQA** Fusion is a common technique applied in many VQA works to integrate language and image features into a joint embedding for answer prediction. There are various types of fusion strategies for multi-modalities including simple concatenation and summation. For example, [22] concatenated question and object features together and pass the joint vectors to a bidirectional GRU for further processes. However, the recent bilinear fusion methods are more effective at capturing higher level of interactions between different modalities and have less parameters in calculation. For example, based on the tensor decomposition proposed in [3], [4] proposed a block-term decomposition of the projection tensor in bilinear fusion. [5] applied this block-term fusion in their proposed MuRel networks, where sequences of MuRel cells are stacked together to fuse visual features and question features together.

**Computational Attention-based VQA** Apart from fusion techniques, attention mechanisms are also commonly applied in VQA for the integration of multi-modal features. Such attention mechanisms include soft attention mechanism like [11,1] using softmax to generate attention weights over object regions and question words, self attention mechanism like [24,18] that applied dot products on features of each mode, and co-attention mechanisms like in [16,6] using linguistic features to guide attentions of visual features or vice versa.

**Graph Representations in VQA** In recent years, more works have been proposed to integrate graph representations of images in VQA model. [19] proposed a question specific graph-based model where objects are identified and connected with each other if their relationships are implied in the given question. There are also works use scene graph in VQA like we did. [21] integrates scene graphs together with functional programs for explainable reasoning steps. [8] claimed only partial image scene graphs are effective for answer prediction and proposed a selective system to choose the most important path in a scene graph and use the most probable destination node features to predict an answer. However, these works did not apply their models on GQA.

**REXUP's Contribution** In this work, we move away from the classical attention and traditional fusion network, which have been widely used in simple photo-realistic VQA tasks and focus mainly on the interaction between visual pixel features from an image and question embeddings. Instead, we focus on proposing a deeper reasoning solution in visual-and-language analysis, as well as complex object-relationship detection in complex photo-realistic images. We propose a new deep reasoning VQA model that can be worked well on complex images by processing both image objects features and scene graph features and integrating those with super-diagonal fusion compositional attention networks.

## 3   Methodology

The REXUP network contains two parallel branches, object-oriented branch and scene-graph oriented branch, shown in Fig. 1a. Each branch contains a sequence of REXUP cells where each cell operates for one reasoning step for the answer
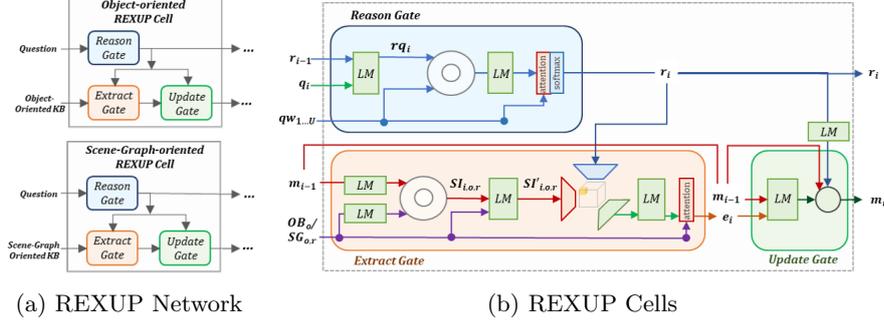
(a) REXUP Network                    (b) REXUP Cells

Fig. 1: **REXUP Network and REXUP cell.** (a) The REXUP network includes two parallel branches, object-oriented *(top)* and scene graph-oriented *(bottom)*. (b) A REXUP cell contains reason, extract, and update gate which conduct multiple compositional reasoning and super-diagonal fusion process

prediction. As shown in Fig. 1b, each REXUP cell includes a reason, an extract and an update gate. At each reasoning step, the reason gate identifies the most important words in the question and generates a current reasoning state with distributed attention weights over each word in the question. This reasoning state is fed into the extract gate and guides to capture the important objects in the knowledge base, retrieving information that contains the distributed attention weights over objects in the knowledge base. The update gate takes the reasoning state and information from extract gate to generate the current memory state.

### 3.1    Input Representation

Both Object-oriented branch and Scene graph-oriented branch take question and knowledge base as inputs; image object-oriented knowledge base (OKB) and scene-graph-oriented knowledge base (SGKB). For a question $q \in Q$ with maximum $U$ words, contextual words are encoded via a pre-trained $300d$ Glove embedding and passed into bi-LSTM to generate a sequence of hidden states $qw_{1...U}$ with $d$ dimension for question contextual words representation. The question is encoded by the concatenation of the last backward and forward hidden states, $\overleftarrow{qw_1}$ and $\overrightarrow{qw_U}$. Object features are extracted from a pre-trained Fast-RCNN model, each image contains at most 100 regions represented by a $2048d$ object feature. For each $o_{th}$ object in an image, linear transformation converts the object features with its corresponding coordinates to a $512d$ object region embedding. The SGKB is the matrix of scene graph objects each of which is in 900 dimensions after concatenating with their corresponding attribute and relation features. To encode the scene graph object features, all the objects names, their attributes and relations in the scene graph are initialized as $300d$ Glove embedding. For each object's attributes, we take the average of those attributes features $A$. For each object's relations, we first average each relation feature $r_s \in R$ and the subject feature $o_j \in O$ that it is linked to, and then average all

such relation-subject features that this object $o_n \in O$ has as its final relation feature. We concatenate the object feature, attribute feature and relation feature together as one scene graph object feature $SG_{o,r}$ of the whole scene graph.

### 3.2 REXUP Cell

With the processed input, each branch consists of a sequence of REXUP cells where each cell operates for one reasoning step for the answer prediction.

**Reason Gate** At each reasoning step, the reason gate in each REXUP cell $i = 1, ..., P$ takes the question feature $q$, the sequence of question words $qw_1, qw_2, ..., qw_U$ and the previous reasoning state $r_{i-1}$ as inputs. Before being passed to the reason gate, each question $q$ is processed through a linear transformation $q_i = W_i^{d \times 2d} q + b_i^d$ to encode the positional-aware question embedding $q_i$ with $d$ dimension in the current cell. A linear transformation is then processed on the concatenation of $q_i$ and the previous reasoning state $r_{i-1}$,

$$rq_i = W^{d \times 2d} [r_{i-1}, q_i] + b^d \tag{1}$$

in order to integrate the attended information at the previous reasoning step into the question embedding at the current reasoning step.

Then an element-wise multiplication between $rq_i$ and each question word $qw_u$, where $u = 1, 2, ..., U$, is conducted to transfer the information in previous reasoning state into each question word at the current reasoning step, the result of which will be processed through a linear transformation, yielding a sequence of new question word representations $ra_{i,1}, ..., ra_{i,u}$ containing the information obtained in previous reasoning state. A softmax is then applied to yield the distribution of attention scores $rv_{i,1}, ..., rv_{i,u}$ over question words $qw_1, ..., qw_u$.

$$ra_{i,u} = W^{1 \times d}(rq_i \odot qw_u) + b \tag{2}$$

$$rv_{i,u} = softmax(ra_{i,u}) \tag{3}$$

$$r_i = \sum_{u=1}^{U} rv_{i,u} \cdot qw_u \tag{4}$$

The multiplications of each $rv_{i,u}$ and question word $qw_u$ are summed together and generates the current reasoning state $r_i$ that implies the attended information of a question at current reasoning step.

**Extract Gate** The extract gate takes the current reasoning state $r_i$, previous memory state $m_{i-1}$ and the knowledge base features as inputs. For the OKB branch, knowledge base features are the object region features $OB_o$, and for the SGKB branch, knowledge base features are the scene graph features $SG_{o,r}$. For each object in the knowledge base, we first multiplied its feature representation with the previous memory state to integrate the memorized information at the

previous reasoning step into the knowledge base at the current reasoning step, the result of which is then concatenated with the input knowledge base features and projected into $d$ dimensions by a linear transformation. This interaction $SI'_{i,o,r}$ generates the knowledge base features that contains the attended information memorized at the previous reasoning step as well as the yet unattended information of knowledge base at current reasoning step. The process of the extract gate in the SGKB branch can be shown in the following equations, where the interaction $SI'_{i,o,r}$ contains the semantic information extracted from the object-oriented scene graph.

$$SI_{i,o,r} = \left[ W_m^{d \times d} m_{i-1} + b_m^d \right] \odot \left[ W_s^{d \times d} SG_{o,r} + b_s^d \right] \tag{5}$$

We then make $SI'_{i,o,r}$ interact with $r_i$ to let the attended question words guide the extract gate to detect important objects of knowledge base at the current reasoning step. In the SGKB branch, such integration is completed through a simple multiplication as shown in (7).

$$SI'_{i,o,r} = W^{d \times 2d} [SI_{i,o,r}, SG_{o,r}] + b^d \tag{6}$$

$$ea_{i,o,r} = W^{d \times d} (r_i \odot SI'_{i,o,r}) + b^d \tag{7}$$

However, in OKB branch, $SG_{o,r}$ in Equation (5) and (6) is replaced with the object region features $OB_o$, and generated interaction $I'_{i,o}$, which will be integrated with $r_i$ through a super-diagonal fusion [4] as stated in Equation (8), where $\theta$ is a parameter to be trained. Super-diagonal fusion projects two vectors into one vector with $d$ dimension through a projection tensor that would be decomposed into three different matrices during calculation in order to decrease the computational costs while boosting a stronger interaction between input vectors. The resulted $F_{r_i, I'_{i,o}}$ is passed via a linear transformation to generate $ea_{i,o}$.

$$F_{r_i, I'_{i,o}} = SD(r_i, I'_{i,o}; \theta) \quad \text{and} \quad ea_{i,o} = W^{d \times d} F_{r_i, I'_{i,o}} + b^d \tag{8 and 9}$$

Similar to the process in the reason gate, $ea_{i,o,r}$ and $ea_{i,o}$ are then processed by softmax to get the distribution of attention weights for each object in the knowledge base. The multiplications of each $ea_{i,o,r}/ea_{i,o}$ and knowledge base $SG_{o,r}/OB_o$ are summed together to yield the extracted information $e_i$.

$$ev_{i,o,r} = softmax(ea_{i,o,r}) \quad \text{and} \quad ev_{i,o} = softmax(ea_{i,o}) \tag{10}$$

$$e_i = \sum_{o=1}^{O} ev_{i,o,r} \cdot SG_{o,r} \quad \text{and} \quad e_i = \sum_{o=1}^{O} ev_{i,o} \cdot OB_o \tag{11}$$

**Update Gate** We apply a linear transformation to the concatenation of the extracted information $e_i$ and previous memory state $m_{i-1}$ to get $m_i^{prev}$.

$$m_i^{prev} = W^{d \times 2d} [e_i, m_{i-1}] + b^d \tag{12}$$

$$m_i = \sigma(r'_i)m_{i-1} + (1 - \sigma(r'_i))m'_i \tag{13}$$

To reduce redundant reasoning steps for short questions, we applied sigmoid function upon $m_i^{prev}$ and $r'_i$, where $r'_i = W^{1 \times d}r_i + b^1$, to generate the final memory state $m_i$.

The final memory states generated in the OKB branch and SGKB branch respectively are concatenated together as the ultimate memory state $m_P$ for overall $P$ reasoning steps. $m_P$ is then integrated with the question sentence embedding $q$ for answer prediction. In this work, we set $P = 4$.

## 4 Evaluation

### 4.1 Evaluation Setup

**Dataset** Our main research aim is proposing a new VQA model that provides not only complex object-relationship detection capability, but also deep reasoning ability. Hence, we chose the GQA that covers 1) complex object-relationship: 113,018 photo-realistic images and 22,669,678 questions of five different types, including *Choose, Logical, Compare, Verify and Query*, and 2) deep reasoning tasks: over 85% of questions with 2 or 3 reasoning steps and 8% of questions with 4+ reasoning steps. The GQA is also annotated with scene graphs extracted from the Visual Genome [15] and functional programs that specify reasoning operations for each pair of image and question. The dataset is split into 70% training, 10% validation, 10% test-dev and 10% test set.

**Training Details** The model is trained on GQA training set for 25 epochs using a 24 GB NVIDIA TITAN RTX GPU with 10.2 CUDA toolkit. The average per-epoch training times and total training times are 7377.31 seconds and 51.23 hours respectively. We set the batch size to 128 and used an Adam optimizer with an initial learning rate of 0.0003.

### 4.2 Performance Comparison

In Table. 1, we compare our model to the state-of-the-art models on the validation and test-dev sets of GQA. Since the GQA test-dev set does not provide pre-annotated scene graphs, we used the method proposed in [26] to predict relationships between objects and generate scene graphs from images of GQA test-dev set for the evaluation procedure. However, the quality of the generated scene graphs are not as good as the pre-annotated scene graphs in the GQA validation set, which lead to the decreased performance on test-dev. Nevertheless, our model still achieves the state-of-the-art performance with 92.7% on validation and 73.1% on test-dev. Compared to [1,11,23] that only used the integration between visual pixel features and question embedding through attention mechanism, our model applies super-diagonal fusion for a stronger interaction and also integrates the scene graph features with question embedding, which help

to yield much higher performance. Moreover, our model greatly improves over [10], which used the graph representation of objects but concatenated the object features with contextual relational features of objects as the visual features to be integrated with question features through the soft attention. The significant improvement over [10] indicates that the parallel training of OKB and SGKB branch can successfully capture the structural relationships of objects in images.

Table 1: State-of-the-art performance comparison on the GQA dataset

| Methods | Val | Test-dev |
|---|---|---|
| CNN+LSTM [12] | 49.2 | - |
| Bottom-Up [1] | 52.2 | - |
| MAC [11] | 57.5 | - |
| LXMERT [23] | 59.8 | 60.0 |
| single-hop [10] | 62 | 53.8 |
| single-hop+LCGN [10] | 63.9 | 55.8 |
| **Our Model** | **92.7** | **73.1** |

Table 2: Results of ablation study on **validation** and **test-dev** set of GQA. 'O' and 'X' refers to the existence and absence of scene-graph oriented knowledge branch($SGKB$) and super-diagonal($SD$) fusion applied in object-oriented knowledge branch($OKB$) branch respectively

| # | OKB | SD | SGKB | Val | Test-dev |
|---|---|---|---|---|---|
| 1 | O | X | X | 62.35 | 56.92 |
| 2 | O | O | X | 63.10 | 57.25 |
| 3 | O | X | O | 90.14 | 72.38 |
| **4** | **O** | **O** | **O** | **92.75** | **73.18** |

### 4.3   Ablation study

We conducted the ablation study to examine the contribution of each component in our model. As shown in Table 2, integrating object-oriented scene graph features is critical in achieving a better performance on the GQA. Using only OKB branch leads to a significant drop of 29.65% in the validation accuracy and 15.93% in the test-dev accuracy. The significant performance decrease also proves the importance of semantic information of objects' structural relationships in VQA tasks. Moreover, applying the super-diagonal fusion is another key reason of our model's good performance on GQA. We compared performances of models that apply super-diagonal fusion and models that apply element-wise multiplication. The results show that using element-wise multiplication causes a drop of 2.61% on the validation set and 0.8% on the test-dev set. It still shows that the concrete interaction between image features and question features generated by super-diagonal fusion contributes to an improved performance on the GQA.

### 4.4   Parameter Comparison

Sequences of REXUP cells will lead to sequential reasoning steps for the final answer prediction. The three gates in each cell are designed to follow questions'

Table 3: Parameter Testing with different number of the REXUP cell

| # of cells | Val | Test-dev |
|:---:|:---:|:---:|
| 1 | 90.97 | 72.08 |
| 2 | 90.98 | 72.13 |
| 3 | 92.67 | 72.56 |
| **4** | **92.75** | **73.18** |

compositional structures and retrieve question-relevant information from knowledge bases at each step. To reach the ultimate answer, a few reasoning steps should be taken, and less cells are insufficient to extract the relevant knowledge base information for accurate answer prediction, especially for compositional questions with longer length. In order to verify this assumption, we have conducted experiments to examine the model's performances with different numbers of REXUP cells in both branches. The results of different performances are shown in Table 3. From the result, we can see that the prediction accuracy on both validation and test-dev set will gradually increase (90.97% to 92.75% on validation and 72.08% to 73.18% on test-dev) as the cell number increases. After experiment, we conclude that four REXUP cells are best both for clear presentation of reasoning capabilities and a good performance on the GQA.

### 4.5   Interpretation

To have a better insight into the reasoning abilities of our model, we extract the linguistic and visual attention weights our model computes at each reasoning step to visualize corresponding reasoning processes. Taking the first row in Fig. 2 as an example, at the first reasoning step, concrete objects - man's hand and head obtain high visual attention score. When it comes to the second and third reasoning step, linguistic attention focuses on *wearing* and corresponding visual attention focuses on man's shirt and pants. This indicates that our model's abilities in capturing the underlying semantic words of questions as well as detecting relevant objects in image for answer prediction. Moreover, our model's good understanding of both images and questions is also shown when given different questions for a same image. For example, in the second row in Fig. 2, the model successfully captures the **phone** in image for the question, but for images of third row in Fig. 2, the **dog** is detected instead. We also found that sometimes our predicted answer is correct even though it's different from the answer in dataset. For example, in the first image of Fig. 3a, our model assigns a high visual attention score to wetsuit in image when question words *person* and *wearing* are attended. Our model then gives the prediction **wetsuit**, which is as correct as **shoe** considering the given image and question. Similarly, in the second image of Fig. 3a, both white bus and red bus are spatially on the right of garbage. Our model captures both buses but assigns more attention to the red bus that is more obvious on the picture and predicts *no*, which is also a correct answer to the question. In addition, we found that in some cases our model's
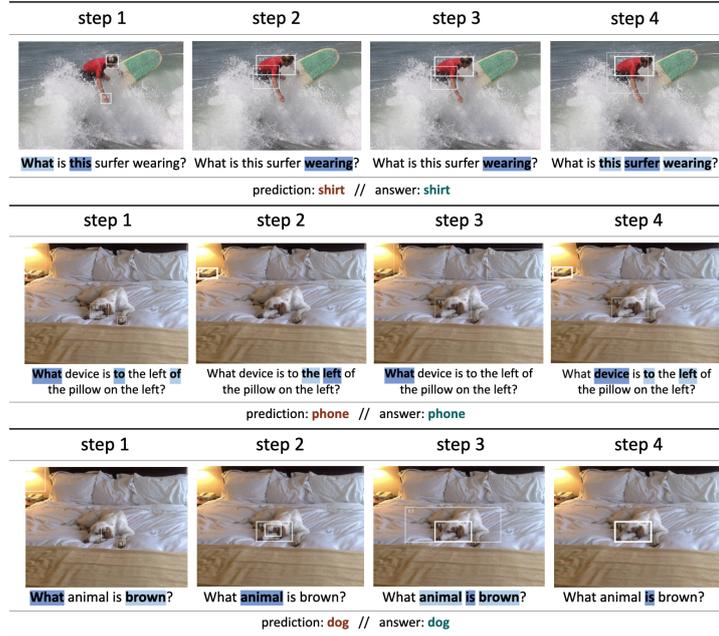
Fig. 2: Visualization of important image objects and question words at each reasoning step. Object regions with high attention weights are framed with white bounding boxes. The thicker the frame, the more important the object region is. Question words with high attention weights are colored blue in the question.

answer is comparatively more accurate than the annotated answer in dataset. For example, for first image of Fig. 3b, *pen*, as a small area surrounded by fence to keep animal inside, is more accurate than the annotated answer *yard*. Likewise, the bed and quilt are actually different in shape but both in white color, which makes our model's answer correct and the ground truth answer incorrect.

## 5    Conclusion

In conclusion, our REXUP network worked well in both capturing step-by-step reasoning process and detecting a complex object-relationship in photo-realistic images. Our proposed model has achieved the state-of-the-art performance on the GQA dataset, which proves the importance of structural and compositional relationships of objects in VQA tasks. Extracting the semantic information of scene graphs and encoding them via textual embeddings are efficient for the model to capture such structural relationships of objects. The parallel training of two branches with object region and scene graph features respectively help the model to develop comprehensive understanding of both images and questions.

(a)                                              (b)

Fig. 3: Figure 3a shows examples when the ground truth answer and our prediction are both correct to the given question. Figure 3b shows examples when our prediction is more accurate than the ground truth answer in dataset

# References

1. Anderson, P., He, X., Buehler, C., Teney, D., Johnson, M., Gould, S., Zhang,L.: Bottom-up and top-down attention for image captioning and visual question answering. In: IEEE conference on computer vision and pattern recognition. pp.6077-6086 (2018)
2. Antol, S., Agrawal, A., Lu, J., Mitchell, M., Batra, D., Lawrence Zitnick, C., Parikh,D.: Vqa: Visual question answering. In: IEEE international conference on computer vision. pp. 2425-2433 (2015)
3. Ben-Younes, H., Cadene, R., Cord, M., Thome, N.: Mutan: Multimodal tucker fusion for visual question answering. In: IEEE international conference on computer vision. pp. 2612-2620 (2017)
4. Ben-Younes, H., Cadene, R., Thome, N., Cord, M.: Block: Bilinear super-diagonal fusion for visual question answering and visual relationship detection. In: AAAI Conference on Artificial Intelligence. vol. 33, pp. 8102-8109 (2019)
5. Cadene, R., Ben-Younes, H., Cord, M., Thome, N.: Murel: Multimodal relational reasoning for visual question answering. In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 1989-1998 (2019)
6. Gao, P., Jiang, Z., You, H., Lu, P., Hoi, S.C.H., Wang, X., Li, H.: Dynamic fusion with intra- and inter-modality attention flow for visual question answering. In:IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2019)
7. Goyal, Y., Khot, T., Summers-Stay, D., Batra, D., Parikh, D.: Making the v in vqa matter: Elevating the role of image understanding in visual question answering.In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 6904-6913(2017)
8. Haurilet, M., Roitberg, A., Stiefelhagen, R.: Its not about the journey; its about the destination: Following soft paths under question-guidance for visual reasoning.In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 1930-1939(2019)
9. Hu, R., Andreas, J., Darrell, T., Saenko, K.: Explainable neural computation via stack neural module networks. In: European conference on computer vision(ECCV). pp. 53-69 (2018)
10. Hu, R., Rohrbach, A., Darrell, T., Saenko, K.: Language-conditioned graph networks for relational reasoning. In: IEEE International Conference on Computer Vision. pp. 10294-10303 (2019)
11. Hudson, D.A., Manning, C.D.: Compositional attention networks for machine rea-soning. In: International Conference on Learning Representations (2018)

12. Hudson, D.A., Manning, C.D.: Gqa: A new dataset for real-world visual reasoning and compositional question answering. In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 6700-6709 (2019)
13. Johnson, J., Hariharan, B., van der Maaten, L., Fei-Fei, L., Lawrence Zitnick,C., Girshick, R.: Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 2901-2910 (2017)
14. Kim, J.H., On, K.W., Lim, W., Kim, J., Ha, J.W., Zhang, B.T.: Hadamard product for low-rank bilinear pooling. In: International Conference on Learning Representations (2016)
15. Krishna, R., Zhu, Y., Groth, O., Johnson, J., Hata, K., Kravitz, J., Chen, S.,Kalantidis, Y., Li, L.J., Shamma, D.A., et al.: Visual genome: Connecting language and vision using crowd sourced dense image annotations. International Journal of Computer Vision123(1), 32-73 (2017)
16. Liu, F., Liu, J., Fang, Z., Hong, R., Lu, H.: Densely connected attention flow for visual question answering. In: 28th International Joint Conference on Artificial Intelligence. pp. 869-875 (2019)
17. Lu, J., Yang, J., Batra, D., Parikh, D.: Hierarchical question-image co-attention for visual question answering. In: Advances In Neural Information Processing Systems.pp. 289-297 (2016)
18. Nguyen, D.K., Okatani, T.: Improved fusion of visual and language representations by dense symmetric co-attention for visual question answering. In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 6087-6096 (2018)
19. Norcliffe-Brown, W., Vafeias, S., Parisot, S.: Learning conditioned graph structures for interpretable visual question answering. In: Advances in Neural Information Processing Systems. pp. 8334-8343 (2018)
20. Perez, E., Strub, F., De Vries, H., Dumoulin, V., Courville, A.: Film: Visual reasoning with a general conditioning layer. In: Thirty-Second AAAI Conference on Artificial Intelligence (2018)
21. Shi, J., Zhang, H., Li, J.: Explainable and explicit visual reasoning over scene graphs. In: IEEE Conference on Computer Vision and Pattern Recognition. pp.8376-8384 (2019)
22. Shrestha, R., Kafle, K., Kanan, C.: Answer them all! toward universal visual question answering models. In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 10472-10481 (2019)
23. Tan, H., Bansal, M.: Lxmert: Learning cross-modality encoder representations from transformers. In: 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing(EMNLP-IJCNLP). pp. 5103-5114 (2019)
24. Yu, Z., Yu, J., Cui, Y., Tao, D., Tian, Q.: Deep modular co-attention networks for visual question answering. In: IEEE Conference on Computer Vision and Pattern Recognition. pp. 6281-6290 (2019)
25. Yu, Z., Yu, J., Fan, J., Tao, D.: Multi-modal factorized bilinear pooling with co-attention learning for visual question answering. In: IEEE international conference on computer vision. pp. 1821-1830 (2017)
26. Zellers, R., Yatskar, M., Thomson, S., Choi, Y.: Neural motifs: Scene graph parsing-with global context. In: Proceedings of the IEEE Conference on Computer Visionand Pattern Recognition. pp. 5831-5840 (2018)