# FROM SHALLOW TO DEEP: COMPOSITIONAL REASONING OVER GRAPHS FOR VISUAL QUESTION ANSWERING

*Zihao Zhu*[1]

[1]University of Chinese Academy of Sciences

## ABSTRACT

In order to achieve a general visual question answering (VQA) system, it is essential to learn to answer deeper questions that require compositional reasoning on the image and external knowledge. Meanwhile, the reasoning process should be explicit and explainable to understand the working mechanism of the model. It is effortless for human but challenging for machines. In this paper, we propose a Hierarchical Graph Neural Module Network (HGNMN) that reasons over multi-layer graphs with neural modules to address the above issues. Specifically, we first encode the image by multi-layer graphs from the visual, semantic and commonsense views since the clues that support the answer may exist in different modalities. Our model consists of several well-designed neural modules that perform specific functions over graphs, which can be used to conduct multi-step reasoning within and between different graphs. Compared to existing modular networks, we extend visual reasoning from one graph to more graphs. We can explicitly trace the reasoning process according to module weights and graph attentions. Experiments show that our model not only achieves state-of-the-art performance on the CRIC dataset but also obtains explicit and explainable reasoning procedures.

***Index Terms***— visual question answering, graph neural modules, compositional reasoning, multi-layer graphs

## 1. INTRODUCTION

One of the goals of AI is to learn to "see" and "talk", which is consistent with Visual Question Answering (VQA) task — aiming to answer natural language questions about an image. Most existing works [1, 2, 3] focused on "shallow" questions which are answerable by solely referring to the visible content of the image. For example, question Q1 in Fig.1 only requires recognizing the color of the helmet to answer, without multi-step reasoning incorporating external knowledge.

However, the "deeper" questions (Q2 in Fig.1) are still obstacles for these methods. To answer this question, a desirable agent should be able to understand the semantic of the question, perceive the visual content (e.g. `helmet`,`boy`), incorporate commonsense knowledge (e.g. `<helmet,Used For,protect head>`) and finally compositionally rea-



**Shallow Question:**
**Q1:** Is there a helmet that is blue or green?
**Answer:** No
**Deep Question:**
**Q2:** What can the hat that the batter is wearing be used for?
**Answer:** Protecting head

**Fig. 1**. Examples of shallow and deep visual questions.

sons over these clues to predict the correct answer. Besides, it should present the decision-making process to better understand the model's underlying working mechanism. Therefore, how to extend shallow visual understanding to deeper compositional reasoning and meanwhile provide an explainable diagnosis are essential to achieve a long-standing general VQA goal.

Most of existing VQA models [4, 5, 6, 7, 8, 9] that exhibit reasoning capabilities can be divided into three categories. Firstly, attention-based methods [4, 5] stack attention layer that focuses on image regions relevant to the question. The reasoning process can be post-hoc extracted by observing attention weights. However, it is implicit reasoning process because the intermediate heat map cannot clarify what the current decision step is. Secondly, memory-based methods [6, 7] perform read and write operations to external memory module iteratively. It implements reasoning by modeling interaction between multiple parts of data over several passes with attention mechanism, which still faces shortcomings of implicit reasoning. Thirdly, module-based methods [8, 9] implement reasoning procedure by parsing the question into a layout of sub-tasks that are carried out by separate neural sub-networks. The reasoning steps can be explicitly defined by each module, but relying on strong layout supervision. Moreover, most existing methods rarely utilize commonsense knowledge that is essential for deeper reasoning.

In this work, we propose a Hierarchical Graph Neural Module Network (HGNMN) to address the above problems by compositional reasoning over multi-layer graphs. Compared with holistic features, graph can model the relationships between objects and makes it easier for a model to reasons from one node to another along the edges [10]. Therefore, we first encode an image by multi-layer graphs from the visual, semantic and commonsense views, which contain the
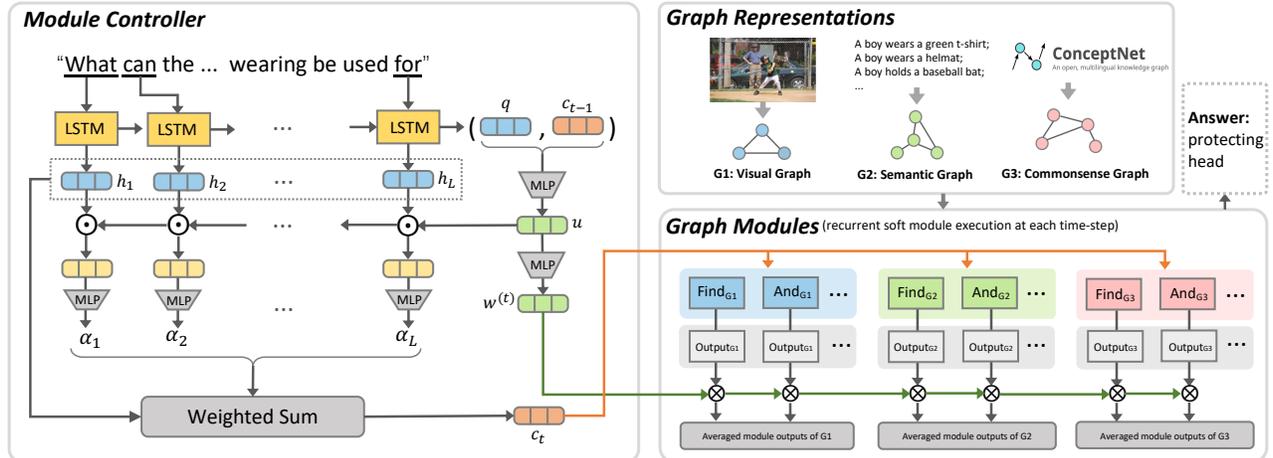
**Fig. 2**. An overview of our model at time-step $t$. It consists of graph representations, graph modules and module controller

clues supporting the answer from different modalities. Second, we decompose the inference process into several subtasks based on a set of well-designed neural modules, which perform specific functions over graphs under the guidance of the semantic information of the question. Third, we propose a module controller that supplies soft module weights and input query at each reasoning step, which makes our model fully differentiable and trainable using gradient descent without resorting to expert layouts. Experimental results on the CRIC dataset [11] demonstrate that HGNMN has comparable performance to existing modular networks and monolithic networks. Meanwhile, the reasoning process can be explicitly traced via module weights and graph attention mechanism.

## 2. PROPOSED APPROACH

In this section, we elaborate on the proposed Hierarchical Graph Neural Module Network (HGNMN) for answering deeper visual questions. Fig. 2 shows the overall architecture. It consists of three components: (1) Graph Representations. Since the clues supporting answer the question may exist in different modalities, we encode the image by constructing multi-layer graphs from the visual, semantic and commonsense views respectively; (2) Graph Modules, a set of well-designed neural modules, that can conduct reasoning over hierarchical graphs; (3) Module Controller guides the execution of each module at each reasoning step.

### 2.1. Graph Representations

Taking the inspiration from recent graph-based methods [6, 12, 13] in visual and language tasks, we first encode the image as multi-layer graphs for unifying the structure representation of different modalities, including visual graph, semantic graph and commonsense graph.

**Visual Graph** represents visual objects and their relationships of the image since most of the questions in VQA are visually related. We construct a fully-connected visual graph over a set of objects identified by Faster-RCNN [14]. Each node corresponds to a detected object and encoded by visual feature. Each edge denotes the relationships between objects, encoded by relative spatial feature.

**Semantic Graph** keeps high-level abstraction of the objects and relationships within natural language also provides essential semantic information. We first generate captions of the image with DenseCap [15]. Then we construct a semantic graph using SPICE [16], where node represents the name or attribute of an object and edge represents the relationship between them. We use the averaged GloVe [17] embeddings to represent nodes and edges.

**Commonsense Graph** contains knowledge related to the image and question that is used for answering open-domain visual questions. We use ConceptNet [18] as original knowledge base, which can be seen as a large set of triples of the form $(h, r, t)$, where $h$ and $t$ represent head and tail entities, $r$ represents relationship between them. We select 10 types of relations to reduce the computation cost following [11]. A desirable knowledge retrieval should include most of the useful information while ignore the irrelevant ones. To achieve this goal, we first use labels of visual objects as keys to extract corresponding entities of ConceptNet. Then we retrieve the first-order subgraph using these selected nodes from ConceptNet, which includes all edges connecting with at least one candidate node, i.e. the label is either $h$ or $t$. However, the subgraph contains much irrelevant information. Each object and extracted triplet is associated with probability scores, denoted as $S_l$ and $S_t$. We assign $a \cdot S_l + b \cdot S_t$ as the final score of the edge, where $a$ and $b$ are hyper-parameters that used to adjust the importance of the $S_l$ and $S_t$. Then we obtain the commonsense graph by ranking and selecting top-K edges along with the nodes according to scores. All nodes

and edges are encoded by averaged GloVe embeddings[17].

## 2.2. Graph Modules

We first define the `Find`, `And`, `Filter` modules to attend objects that are relevant to some subjects, attributes or logic information of the question following [9, 19]. However, deeper questions commonly cover relationships (e.g. geometric positions or semantic interactions between objects) beyond mere object detection, so we propose the `Relate` module to transfer node attention maps along the edge in one graph guided by relation information of the question. Moreover, deeper questions require collecting evidence from different modalities, so it is necessary to associate multi-layer graphs. We additionally propose the `CrossGraph` module to extend the reasoning in one graph to more graphs. The `Describe` module is proposed to obtain the graph features because the attentive node maps need to be transformed to an embedding to predict the answer. Due to some questions may not require complete $T$-step reasoning, we also introduce a `NoOp` module proposed by [9], which can be used to pad reasoning steps to maximum length $T$. We will describe each kind of module in detail below. Modules of the same type only differ in inputs and parameters.

**Find [X, c]** is proposed to attend the relevant objects given the input query and outputs an attention map $\mathbf{a} \in \mathbb{R}^n$ over $n$ graph nodes. We first transform the input query $\mathbf{c}$ and node features $\mathbf{X}$ into the same dimensions and then fuse them together to pass a MLP:

$$\mathbf{a} = \text{softmax}(f_{mlp}(F(W_1\mathbf{X}, W_2\mathbf{c}))), \tag{1}$$

where $W_1, W_2$ are weight matrices (as well as $W_3, \ldots, W_{10}$ mentioned below), $F(\mathbf{x}, \mathbf{y}) = \text{ReLU}(\mathbf{x} + \mathbf{y}) - (\mathbf{x} - \mathbf{y})^2$ is multimodal fusion function proposed in [19].

**And [$\mathbf{a}_1, \mathbf{a}_2$]** aims to combine attention maps $\mathbf{a}_1, \mathbf{a}_2$ generated from previous reasoning steps. We implement it by adding two input attention weights, i.e. $\mathbf{a} = \mathbf{a_1} \oplus \mathbf{a_2}$, where $\oplus$ is element-wise addition.

**Filter [$\mathbf{a}, \mathbf{c}$]** is proposed to find nodes relevant to the input query on the basis of the output from the previous step, which is used to further locate objects according to new input information. We implement it based on `Find` module, i.e. $\mathbf{a} = \mathbf{a} \oplus \text{Find}(\mathbf{c})$, where $\oplus$ is element-wise addition.

**Relate [$\mathbf{a}, \mathbf{E}, \mathbf{c}$].** The deeper questions mostly involve the interaction between objects. Therefore it is essential to transfer the current node to adjacent nodes via attended edge. We first find the relevant edges given the input query. The attention weight $W_{ij} \in \mathbf{W}^{n \times n}$ of edge $e_{ij}$ is computed by:

$$W_{ij} = \text{ReLU}(f_{mlp}(W_3\mathbf{c} \odot W_4\mathbf{e}_{ij})), \tag{2}$$

where $\mathbf{e}_{ij} \in \mathbf{E}$ is edge embedding. Thanks to graph structure, we can transfer the node weights along the attentive relations to update attention weights by matrix multiplication: $\mathbf{a} = \text{norm}(\mathbf{W}^T\mathbf{a})$, where $\text{norm}(\cdot)$ is normalization operation.

**CrossGraph**$_{G_m \to G_n}$ **[$\mathbf{a}_m, \mathbf{a}_n, \mathbf{X}_m, \mathbf{X}_n, \mathbf{c}$].** Reasoning within one graph is not enough to answer deep questions. It is essential to associate nodes of different graphs to extend the reasoning process from one graph to multi-layer graphs. We achieve this goal by computing attention map over nodes of graph $G_n$ guided by the node features of graph $G_m$:

$$\mathbf{a'}_n(i) = \text{softmax}(W_7(\tanh(W_5\mathbf{X}_m^T\mathbf{a}_m + W_6\mathbf{X}_n(i))), \tag{3}$$
$$\mathbf{a}_n = \mathbf{a'}_n + \mathbf{a}_n, \tag{4}$$

where $\mathbf{X}_n(i)$ is node feature of the *i-th* node of $G_n$.

**Describe [a, X]** aims to transform the attentive node features to an embedding that summarize the entire graph by computing a weighted sum of node features guided by the input attention, i.e. $\mathbf{y} = \mathbf{X}^T\mathbf{a}$.

**NoOp [a]** just outputs the input $\mathbf{a}$ without transforming, which is used to pad reasoning steps to a maximum length $T$.

## 2.3. Module Controller

To make the network fully differentiable, we draw inspire from [9] to propose a module controller to make soft layout selection via a soft module weights distribution and supplies input query at each reasoning step $t$. Then we execute all modules of each graph and average the outputs according to the module weights as the input for the next step. We take one of the graphs as an example to describe the procedure in detail below.

Specifically, we first encode the question by LSTM to get all word embeddings $\{\mathbf{h}_l\}_{l=1}^L$ and question embedding $\mathbf{q}$. At each time-step $t$, we generate intermediate embedding $\mathbf{u}$ based on question embedding and input query of previous step, i.e $\mathbf{u} = f_{mlp}([\mathbf{q}; \mathbf{c}_{t-1}])$. The input query should include question information that is more relevant to the current reasoning step. Therefore we generate input query by textual attention over words guided by intermediate embedding:

$$\alpha_l = \text{softmax}(f_{mlp}(\mathbf{u} \odot \mathbf{h}_l)) \tag{5}$$

$$\mathbf{c}_t = \sum_{l=1}^L \alpha_l \cdot \mathbf{h}_l, \tag{6}$$

where $\odot$ is element-wise multiplication. To select which module is more important at the current step, we predict soft module weights which resemble a probability distribution over all the modules using MLP:

$$\mathbf{w}^{(t)} = \text{softmax}(f_{mlp}(\mathbf{u})) \tag{7}$$

We execute all the modules and perform a weighted average of their outputs with respect to the module weights:

$$\mathbf{a}^{(t)} = \sum_{m \in M} w_m^{(t)} \cdot \mathbf{a}_m^{(t)}, \tag{8}$$

where $w_m^{(t)} \in \mathbf{w}^{(t)}$, $M$ is module list of each graph and $\mathbf{a}_m^{(t)}$ is the output attention map of *m-th* module. At the final step, we

| Model | Expert layout | Acc |
|---|---|---|
| Q-Type | No | 25.96 |
| Q-Only | No | 39.40 |
| I-Only | No | 14.18 |
| Q+I | No | 48.47 |
| ButtomUp [20] | No | 52.26 |
| RVC-w/o-KG [11] | Yes | 54.68 |
| RVC-$l_{ans}$ [11] | No | 51.20 |
| RVC [11] | Yes | 58.38 |
| HGNMN (full) | No | **60.32** |

**Table 1**. Results on test set of CRIC dataset.

extract graph features from `Describe` modules of different graphs and fuse them to predict the answer using MLP:

$$Ans = f_{mlp}([W_8\mathbf{y}_{G_1}, W_9\mathbf{y}_{G_2}, W_{10}\mathbf{y}_{G_3}, W_{11}\mathbf{q}]) \quad (9)$$

## 3. EXPERIMENTS

### 3.1. Experimental Settings

**Dataset.** We evaluated our model on the CRIC [11], a large VQA dataset contains 1,303,271 deeper questions than other common datasets, which are randomly split into train (60%), validation (20%) and test (20%). It also provides expert layout annotations for each question.

**Implementation Details.** For each image, we extract 36 objects along with features, labels and probability scores from Faster-RCNN [14]. We select the top-10 captions and set $a = 0.7, b = 0.3$ in the graph construction. We truncate or pad the length of question to 20 words. The max reasoning steps $T$ is set to 12. We adopt cross entropy loss to train the model by Adam optimizer with 0.001 learning rate.

### 3.2. Experimental Results

**Comparison with SOTA Methods.** In this section, we evaluate the performance of following methods on the CRIC: (1) **Q-Type**, **Q-Only**, **I-Only** and **Q+I** are four basic baselines proposed in [11]. (2) **BottomUp [20]** implements soft attention on object regions and combines the attended image features and question features to predict the answer. (3) **RVC [11]** builds upon neural module networks with supervision of expert layout. (4) **RVC-w/o-KG [11]** is a variation of RVC that doesn't use the knowledge graph to answer the question. (5) **RVC-$l_{ans}$ [11]** is a variation of RVC that can be trained without supervision of expert layout.

The results are summarized in Table 1. Our model achieves a new SOTA result. In particular, although RVC is trained with expert layout, our model still outperforms RVC by 1.94%. Compared with RVC-$l_{ans}$ that also does not require expert layout, our model improves by 9.12%, which proves the effectiveness of reasoning on multi-layer graphs.

**Ablation Study.** We conduct ablation studies to further investigate the key components of HGNMN. (1) We remove the visual, semantic and commonsense graph from full model

| Model | Acc |
|---|---|
| w/o VG | 57.63 |
| w/o SG | 59.11 |
| w/o KG | 55.25 |

| Model | Acc |
|---|---|
| w/o And | 59.12 |
| w/o Filter | 58.47 |
| w/o Relate | 57.39 |
| w/o CrossGraph | 57.14 |

**Table 2**. Ablation study of hierarchical graphs   **Table 3**. Ablation study of different modules
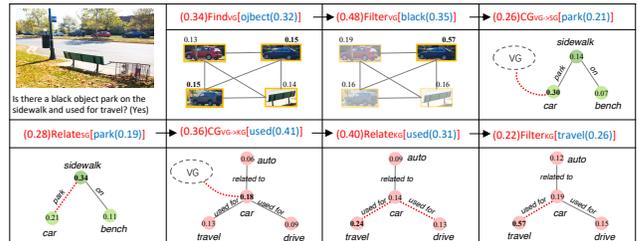


**Fig. 3**. Visualization of intermediate reasoning steps.

respectively to analyze the contributions of each layer of graphs. The results shown in Table 2 all decrease. Thereinto, the commonsense graph is most important since most of the questions rely on external knowledge. (2) We further evaluate the effectiveness of different neural modules in Table 3. We remove the `And`, `Filter`, `Relate` and `CrossGraph` modules respectively. The `Find`, `Describe` and `NoOp` modules are kept because answering any questions at least requires these operations. We find that the `Relate` and `CrossGraph` modules are more important than others, since they can extend reasoning along edges and associate multiple graphs which are essential for multi-step reasoning.

**Qualitative Evaluation** Fig. 3 shows the qualitative evaluation of one example. We visualize the results of the first 7 steps and omit the rest `NoOp` operations. At each reasoning step, we mark the module with the largest weight in red and the most attentive word of the question in blue and present the attention map of corresponding graph. It shows that our model can extract explicit and explainable reasoning process. For example, it locates the "black object" in VG via `Filter` at the 2*nd* step; finds the most relevant object in SG via `CrossGraph` at the 3*rd* step; finds the nodes related to "car" in KG via `Relate` at the 6*th* step.

## 4. CONCLUSION

In this paper, we propose a Hierarchical Graph Neural Module Network (HGNMN) for answering deeper questions of VQA. We encode the image by multi-layer graphs from different views and define a set of graph-based neural modules to extend reasoning from single graph to more graphs. Moreover, it can be trained end-to-end without expert layout supervision. Experimental results on CRIC dataset show that HGNMN outperforms state-of-the-art approaches. Our model is more interpretable by observing its intermediate outputs.

# 5. REFERENCES

[1] Huaizu Jiang, Ishan Misra, Marcus Rohrbach, Erik Learned-Miller, and Xinlei Chen, "In defense of grid features for visual question answering," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10267–10276.

[2] Hedi Ben-Younes, Remi Cadene, Nicolas Thome, and Matthieu Cord, "Block: Bilinear superdiagonal fusion for visual question answering and visual relationship detection," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2019, vol. 33, pp. 8102–8109.

[3] Tuong Do, Thanh-Toan Do, Huy Tran, Erman Tjiputra, and Quang D Tran, "Compact trilinear interaction for visual question answering," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 392–401.

[4] Zichao Yang, Xiaodong He, Jianfeng Gao, Li Deng, and Alex Smola, "Stacked attention networks for image question answering," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 21–29.

[5] Zhou Yu, Jun Yu, Yuhao Cui, Dacheng Tao, and Qi Tian, "Deep modular co-attention networks for visual question answering," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2019, pp. 6281–6290.

[6] Jing Yu, Zihao Zhu, Yujing Wang, Weifeng Zhang, Yue Hu, and Jianlong Tan, "Cross-modal knowledge reasoning for knowledge-based visual question answering," *Pattern Recognition*, vol. 108, pp. 107563, 2020.

[7] Chao Ma, Chunhua Shen, Anthony Dick, Qi Wu, Peng Wang, Anton van den Hengel, and Ian Reid, "Visual question answering with memory-augmented networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 6975–6984.

[8] Wenhu Chen, Zhe Gan, Linjie Li, Yu Cheng, William Wang, and Jingjing Liu, "Meta module network for compositional visual reasoning," *arXiv preprint arXiv:1910.03230*, 2019.

[9] Ronghang Hu, Jacob Andreas, Trevor Darrell, and Kate Saenko, "Explainable neural computation via stack neural module networks," in *Proceedings of the European conference on computer vision*, 2018, pp. 53–69.

[10] Sibei Yang, Guanbin Li, and Yizhou Yu, "Graph-structured referring expression reasoning in the wild," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9952–9961.

[11] Difei Gao, Ruiping Wang, Shiguang Shan, and Xilin Chen, "From two graphs to n questions: A vqa dataset for compositional reasoning on vision and commonsense," *arXiv preprint arXiv:1908.02962*, 2019.

[12] Linjie Li, Zhe Gan, Yu Cheng, and Jingjing Liu, "Relation-aware graph attention network for visual question answering," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 10313–10322.

[13] Zihao Zhu, Jing Yu, Yujing Wang, Yajing Sun, Yue Hu, and Qi Wu, "Mucko: Multi-layer cross-modal knowledge reasoning for fact-based visual question answering," in *Proceedings of the International Joint Conference on Artificial Intelligence*, 2020, pp. 1097–1103.

[14] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Advances in neural information processing systems*, 2015, pp. 91–99.

[15] Justin Johnson, Andrej Karpathy, and Li Fei-Fei, "Densecap: Fully convolutional localization networks for dense captioning," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 4565–4574.

[16] Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould, "Spice: Semantic propositional image caption evaluation," in *European Conference on Computer Vision*. Springer, 2016, pp. 382–398.

[17] Jeffrey Pennington, Richard Socher, and Christopher D Manning, "Glove: Global vectors for word representation," in *Proceedings of the 2014 conference on empirical methods in natural language processing*, 2014, pp. 1532–1543.

[18] Robyn Speer, Joshua Chin, and Catherine Havasi, "Conceptnet 5.5: an open multilingual graph of general knowledge," in *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*, 2017, pp. 4444–4451.

[19] Jiaxin Shi, Hanwang Zhang, and Juanzi Li, "Explainable and explicit visual reasoning over scene graphs," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 8376–8384.

[20] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang, "Bottom-up and top-down attention for image captioning and visual question answering," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 6077–6086.