

Visual Question Answering with Adaptive Question Understanding and Reasoning

Anonymous ICCV submission

Paper ID ****

Abstract

Visual Question Answering (VQA) is of great significance in offering people convenience: one can raise a question for details of objects, or high-level understanding about the scene, over an image. This paper proposes a novel method to address the VQA problem. In contrast to prior works, our method that targets single scene VQA, replies on graph-based techniques and involves reasoning. In a nutshell, our approach is centered on three graphs. The first graph, referred to as inference graph G_I , is constructed via learning over labeled data. The other two graphs, referred to as query graph Q and entity-attribute graph G_{EA} , are generated from natural language query Q_{nl} and image Img , that are issued from users, respectively. As G_{EA} often does not take sufficient information to answer Q , we develop techniques to infer missing information of G_{EA} with G_I . Based on G_{EA} and Q , we provide techniques to find matches of Q in G_{EA} , as the answer of Q_{nl} in Img . Unlike commonly used VQA methods that are based on end-to-end neural networks, our graph-based method shows well-designed reasoning capability, and thus is highly interpretable. We also create a dataset on soccer match (Soccer-VQA) with rich annotations. The experimental results show that our approach outperforms the state-of-the-art method and has high potential for future investigation.

1. Introduction

In recent years, visual question answering (VQA) has received significant attention [15, 19, 8] as it involves multi-disciplinary research, e.g. natural language understanding, visual information retrieving and multi-modal reasoning. The task of VQA is to find an answer to a question Q_{nl} based on the content of an image. There are a variety of applications of VQA, e.g. surveillance video understanding, visual commentator robot, etc. Solving VQA problems usually requires high level reasoning from the content of an image.

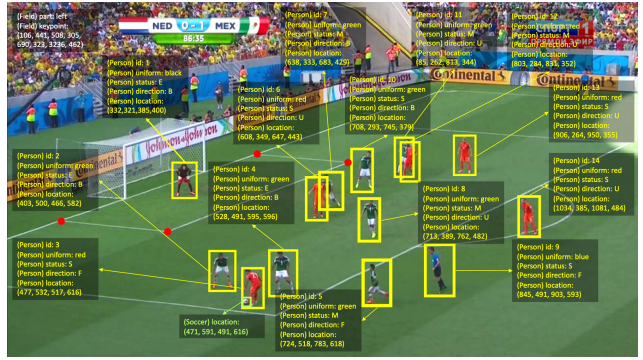


Figure 1: The image is about soccer match, where each person object is associated with attributes: id, uniform color, status (Standing, Moving, Expansion), direction (Backing, Facing, N/A), as well as location, and the soccer object is attributed with location.

Example 1: Figure 1 depicts an image about a soccer match, where two teams are distinguished by red and green uniforms, and each object is associated with a set of attributes. A typical query may ask “How many players are there in the image?”. Though simple, it is a challenging task to efficiently answer the query, since (1) it is often very inefficient to detect all the objects in the given image, following traditional way, while only query related objects are in demand; (2) one not only needs to identify all the *person* objects, but also have to infer their hidden attribute “role”, which is often ambitious.

Motivated by these, one may analyze questions first and detects only those objects as well as their attributes that are in connection with questions, then infer missing values of hidden attributes, and answer questions. In this way, not only query accuracy but also query efficiency are expected to be guaranteed. □

This example suggests that we leverage adaptive query understanding and reasoning to address the VQA problem. While to do this, two critical questions have to be answered. (1) How to understand queries and carry query-related visual tasks? (2) How to infer crucial information to assist

query answering?

Contributions. In contrast to a majority of deep learning based VQA techniques, which not only overlooks correlation between queries and images, but also lacks of necessary reasoning, we propose a novel technique that integrates adaptive query understanding and reasoning. The main contributions of the paper are as follow.

(1) We model images and queries as graphs, and propose to answer visual queries with graph matching. This new representation and answering scheme constitute the base of our techniques.

(2) We introduce a method to guide visual tasks based on reinforcement learning. That is, given an image and a question, our method can identify a set of visual tasks that are question related, and direct subsequent visual processing in a more efficient manner.

(3) We propose approaches to answering visual queries based on reasoning and graph matching. More specifically, we first transform a given image into an entity-attribute graph; we then develop method to infer missing value for question answering; we finally show how to answer queries with graph matching.

(4) We conduct extensive experimental studies to verify the performance of our method. We find that X, Y, and Z.

2. Related Work

We categorize related work into following three parts.

Visual query answering. Current VQA approaches are mainly based on deep neural works. [28] introduces a spatial attention mechanism similar to the model for image captioning. Instead of computing the attention vector iteratively, [26] obtains a global spatial attention weights vector which is then used to generate a new image embedding. [27] proposed to model the visual attention as a multivariate distribution over a grid-structured conditional random field on image regions, thus multiple regions can be selected at the same time. This attention mechanism is called structured multivariate attention in [27]. There has been many other improvements to the standard deep learning method, e.g. [7] utilized Multimodal Compact Bilinear (MCB) pooling to efficiently and expressively combine multimodal features. Another interesting idea is the implementation of Neural Module Networks [1, 10], which decomposes queries into their linguistic substructures, and uses these structures to dynamically instantiate module networks. [22] proposed to build graph over scene objects and question words. The visual graph is similar to ours, but the query graph differs. Note that the method [22] proposed is still a neural network based method as the structured representations are fed into a recurrent network to form the final embedding and the answer is again inferred by a classifier.

Environment Exploration in Visual Field. Reinforcement driven information acquisition is wildly applied in traditional vision domain, like visual object detection [16], face detection [9] and image classification [17]. In visual and language domain, relevant work like [21] achieves image captioning with Embedding Reward. Reinforcement learning preserves ability to effectively select preferred actions, which benefits the system decomposing the problem into a few sub-tasks.

Graph-based visual understanding. [23] proposes a framework to understand events and answer user queries, where underlying knowledge is represented by a spatial-temporal-causal And-Or graph (S/T/C-AOG).

3. Overview of the Approach

We start from representations of images and questions, followed by the overview of our approach.

3.1. Representation of Images and Questions

We use the same representations as [25]. To make the paper self-contained, we cite them as follows (rephrased).

Entity-Attribute Graphs. Entities are typically defined as objects or concepts that exist in the real world. An entity often carries attributes, that describe features of the entity.

Assume a set \mathcal{E} of entities, a set \mathcal{D} of values, a set \mathcal{P} of predicates indicating attributes of entities and a set Θ of types. Each entity e in \mathcal{E} has a *unique ID* and a *type* in Θ .

An *entity-attribute graph*, denoted as EAG, is a set of triples $t = (s, p, o)$, where *subject* s is an entity in \mathcal{E} , p is a *predicate* in \mathcal{P} , and *object* o is either an entity in \mathcal{E} or a value d in \mathcal{D} . It can be represented as a directed edge-labeled graph $G_{EA} = (V, E)$, such that (a) V is the set of nodes consisting of s and o for each triple $t = (s, p, o)$; and (b) there is an edge in E from s to o labeled by p for each triple $t = (s, p, o)$.

An image can be represented as an EAG with detected objects along with their detected attributes, and relationships among objects. This can be achieved via a few visual tasks. While EAG generated directly after image processing is often incomplete, i.e. it may miss some crucial information to answer queries. We hence refer to *entity-attribute graphs* with incomplete information as *incomplete entity-attribute graphs*, and associate nodes with white rectangles, to indicate the missing value of an entity or attribute in EAG. Figure 1(b) is an *incomplete entity-attribute graph*, in which square nodes representing person roles are associated with white rectangle.

Query Graphs. A query graph $Q(u_o)$ is a set of triples (s_Q, p_Q, o_Q) , where s_Q is either a variable z or a function $f(z)$ taking z as parameter, o_Q is one of a value d or z or $f(z)$, and p_Q is a predicate in \mathcal{P} . Here function $f(z)$ is

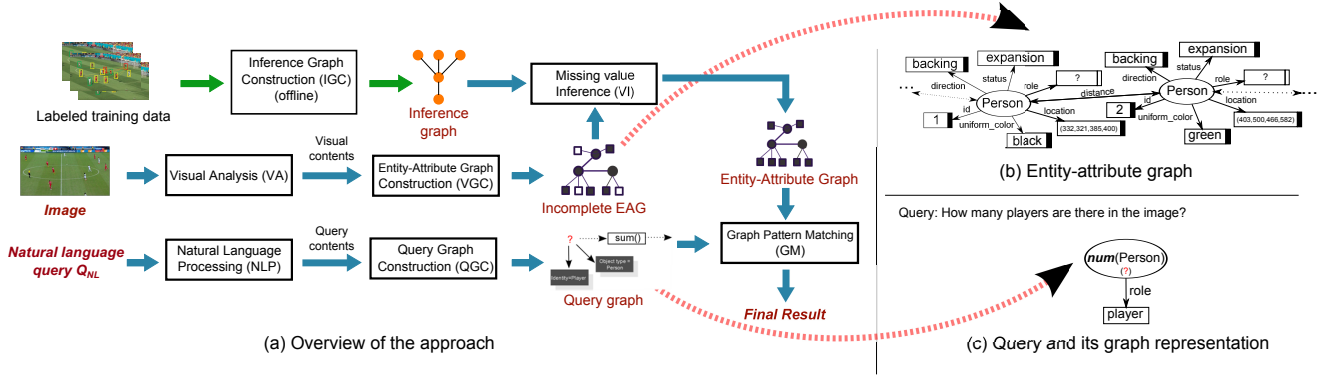


Figure 2: Overview of our approach, and graph-based representation of images and questions

defined by users, and variable z has one of three forms: (a) *entity variable* y , to map to an entity, (b) *value variable* y^* , to map to a value, and (c) *wildcard* $-y$, to map to an entity. Here s_Q can be either y or $-y$, while o_Q can be y , y^* or $-y$. Entity variables and wildcard carry a *type*, denoting the type of entities they represent.

A query graph can also be represented as a graph such that two variables are represented as the same node if they have the same name of y , y^* or $-y$; similarly for functions $f(z)$ and values d . We assume *w.l.o.g.* that $Q(x)$ is connected, *i.e.* there exists an undirected path between u_o and each node in $Q(u_o)$. In particular, u_o is a designated node in $Q(u_o)$, denoting the query focus and labeled by “?”. Take Fig. 2(c) as example. It depicts a query graph that is generated from query “How many players are there in the image?”. Note that the “query focus” u_o carries a function $\text{num}()$ that calculates the total number of *person* entities with *role* “player”.

3.2. Approach Overview

Figure 2(a) presents the overview of our approach. In a nutshell, our approach copes with two types of computations, *i.e.* online computation, which responses to user’s questions and offline computation that trains inference graphs (*i.e.* a classifier) for missing value inference.

For online computation, our approach leverages five modules, *i.e.*

Specifically, our approach consists of two types of

As can be seen, our approach revolves around three graphs: entity-attribute graph, query graph and inference graph. The generation of entity-attribute graph G_{EA} follows three steps. Module VA conducts the first step, *i.e.* image processing, and outputs all the detected objects along with their attributes. Using visual contents produced in step one, module VGC constructs an *incomplete* EAG. In the last step, module VI takes inference graph and *incomplete* EAG as inputs, infer missing information with G_I , and outputs an updated EAG for query answering. The inference graph G_I is used to infer missing values of an *incomplete* EAG. and

constructed by module IGC over training data. As is query-independent, G_I is constructed offline, which warrants the efficiency of our approach. As the other part of input, natural language query Q_{nl} needs to be structured for query evaluation. To this end, Q_{NL} is first parsed via our NLP module, and then structured by module QGC. After $Q(u_o)$ and G_{EA} are generated, our approach employs module GM for matching computation, and returns final result.

It takes an image and a natural language question Q_{NL} as input, and works as following. (1)

As some modules employ existing techniques, to emphasize our novelty, we will elaborate modules VA and VGC in Section ??, modules IGC and VI in Section ??, and module GM in Section ?? with more details.

4. Query Oriented Visual Tasks

In this section, we introduce how we do visual tasks that are in connection with questions.

4.1. Visual Processing

In our approach, we build a structure which selects sub-tasks to form a policy which is based on queries. For instance, with the input *which is the defending team?*, the system first predicts the corresponding visual action sequence, which are *Human Module*, *Gesture Module*, *Direction Module*, *Soccer Module*, *Color Blob Module*, *Field Part Module* and *Graph Indicator*. Guided by such sequence, the image features are then extracted by operating relevant vision tasks. An overview is shown in Figure 3.

4.1.1 Multi-layer LSTM with Attention

The task here is to predict the most suitable action modules sequence a by given questions Q and preference pre . We form the problem of seeking effective answering strategy of question Q and preference pre as a sequence-to-sequence learning problem with attention mechanism. Inspired by [3], we input word feature of questions w_i^q , $i \in ||Q||$ into a LSTM network which is regarded as an encoder


$$e_{ij} = a(s_{j-1}, h_i) \quad (3)$$
$$s_t = f(s_{t-1}, \mathbf{a}_{t-1}, c_i) \quad (6)$$

4

4.1.3 Time and Accuracy Term

To better balance the accuracy and inference time for a given application, we proposed time and accuracy terms in loss function during training process.

$$L_{\tau\alpha}(\theta) = \ell(\theta, A|I, Q) + \gamma \sum_{i \in \|A\|} \alpha(a_i) + (1-\gamma) \sum_{i \in \|A\|} \tau(a_i) \quad (7)$$

where $\tau(\cdot)$ and $\alpha(\cdot)$ represent the pre-tested inference time and inference accuracy, action module sequence A samples from joint distribution $p(A|Q)$, and here $\ell(\cdot)$ is the softmax loss over the predict score. For the preference term γ , it ranges from 0 to 1, which represents the preference over time and accuracy.

4.1.4 Monte Carlo Methods

The task now becomes a policy learning problem. Given a question and preference, output a policy containing a sequence of actions $[a_1, a_2, a_3, \dots, a_n]$. There is no ground truth for each steps, but only a final reward indicates that whether the prediction result is correct based on current policy. We involve the concept of Monte Carlo Methods to learn the policy which guides the vision tasks, and such policy network requires an extra reward value in loss.

$$L_{policy}(\theta) = \sum_{i \in \|A\|} \log \pi(a_i|Q, \theta) \ell(Q, A) \quad (8)$$

where a_i is the action will take, based on current status. $\pi(\cdot)$ is the policy function that maps status to actions, here, the policy is the probability of outputting next action module a_i based on current status. And $\ell(\cdot)$ here is the softmax loss based on the whole action module sequence $[a_1, a_2, a_3, \dots, a_n]$. Since all actions are discrete, which leads to non-differentiable, and back-propagation cannot be used. Policy gradient [12] is used here during training. The object function now becomes the combination of policy gradient loss $L_{policy}(\theta)$ with the time-accuracy-balanced loss $L_{\tau\alpha}$, and optimize it by backpropagation for $L_{\tau\alpha}$, while policy gradient for $L_{policy}(\theta)$.

4.2. Construction of EAG

After objects that are related to questions are identified, we construct a graph structure, denoted as EAG, along the same line as [25].

Example 2: ADD AN EXAMPLE TO ILLUSTRATE PROGRESS IF NECESSARY! \square

5. Reasoning

According to our observation, an *incomplete* EAG isn't well satisfying of answering the query because of the insufficient attributes. To infer the hidden attributes, an inference

graph is constructed accordingly. we briefly introduce the construction below.

5.1. Construction of Inference Graph

To take advantage of the prior information and increase the generalization ability of the proposed model, our inference graph is constructed using Bayesian network. Mathematically, Bayesian network [6] can be described by a pair $\mathfrak{B} = \langle \mathcal{G}, \Theta_{\mathcal{G}} \rangle$. Here, the notation \mathcal{G} is a directed acyclic graph, of which the i -th vertex corresponds to a random variable X_i , and the edge between two connected vertexes indicates the dependency. Additionally, the second item $\Theta_{\mathcal{G}}$ is a set of parameters used to quantify the dependencies in \mathcal{G} . Denoted by $\text{Pa}(X_i)$ the attributes of the parents of X_i , the parameter of X_i is represented by $\theta_{X_i|\text{Pa}(X_i)} = P_{\mathfrak{B}}(X_i|\text{Pa}(X_i))$. With the notations above, the joint probability distribution of Bayesian network is given by:

$$P_{\mathfrak{B}}(X_1, \dots, X_n) = \prod_{i=1}^n P_{\mathfrak{B}}(X_i|\text{Pa}(X_i)) = \prod_{i=1}^n \theta_{X_i|\text{Pa}(X_i)} \quad (9)$$

In our inference graph, the role of Bayesian network is to predict the object class when given the attributes $\{X_i\}_{i=1}^n$ as input. In the sense of probability, the object class is also a variable [11]. Defined by $X_0 = Y$ the class variable, the network now has one extra vertex X_0 . In order to infer the class attribute, and according to the Bayesian rule, our problem becomes:

$$\begin{aligned} P_{\mathfrak{B}}(Y|X) &= \frac{P_{\mathfrak{B}}(Y)P_{\mathfrak{B}}(X|Y)}{P_{\mathfrak{B}}(X)} \\ &= \frac{\theta_{Y|\text{Pa}(X_0)} \prod_{i=1}^n \theta_{X_i|Y, \text{Pa}(X_i)}}{\sum_{y' \in \mathcal{Y}} \theta_{y'|\text{Pa}(X_0)} \prod_{i=1}^n \theta_{X_i|y', \text{Pa}(X_i)}} \end{aligned} \quad (10)$$

where \mathcal{Y} is the set of classes.

5.2. Learning the Structure of Inference Graph

In the context of Naïve Bayes, the structure of $P_{\mathfrak{B}}(Y|X)$ is simplified by taking the class variable as the root, and all attributes are conditionally independent when taking the class as a condition [18]. As a consequence, the attribute class can be explicitly inferred by:

$$P_{\mathfrak{B}}(Y|X) = c \cdot \theta_Y \prod_{i=1}^n \theta_{X_i|Y} \quad (11)$$

where c is a scale factor that makes the calculation being a distribution: $c = \sum_{y' \in \mathcal{Y}} \theta_{y'} \prod_{i=1}^n \theta_{X_i|y'}$.

Note from Eq.(11) that Naïve Bayes simplifies the complexity of Bayesian network. As can be validated by the

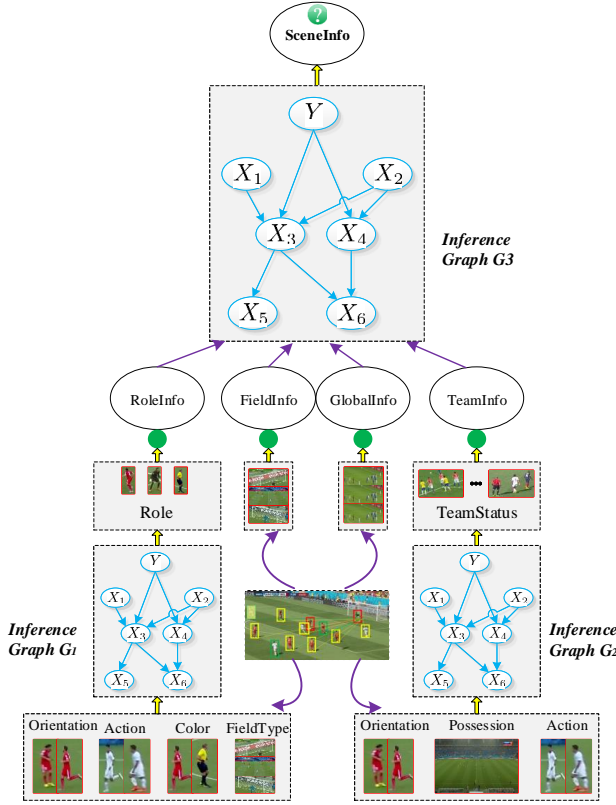


Figure 5: Schematic diagram of inference graph.

experimental results, the simple model works excellently to our problem.

Fig. 5 summarizes the processes of our inference graph, where three graphs are constructed according to the tasks involved. First, the role of a candidate is inferred, in which four different kinds of features are extracted from the scene image. Then, the team status is inferred through the second inference graph, but with different features as input. Next, we use the inferred information, along with the other information can be directly detected from the scene image, to infer the information of the whole scene. The scene information is then fed into the incomplete EAG so that a complete EAG can be obtained.

6. Experimental Studies

In this section, we conducted two sets of experiments to evaluate (1) the performance of our visual processing module, (2) the accuracy of our inference module, and (3) the overall performance of our approach.

Experimental Setting.

DataSet. We used two datasets: (1) Soccer dataset that we annotated; and (2) X dataset from []. We extracted images with subjects of golf and tennis (report Statistics about the dataset). We split Soccer (resp. X) data into two parts: *I* (one third) and *II* (two thirds), and used *II* as training data, and *I*

as testing data.

Queries. We used two sets of questions: (1) the set of questions given in Table 2 for Soccer dataset; and (2) another set of questions listed in Table ?? for X dataset.

(We enlarge the training question scale from 7 into 28, so learning correlation between question and answer does not work at this time. For question details, please refer questionset.txt.rtf.)

Id	Question	Difficulty
Q_{nl1}	Who is holding the soccer?	Easy
Q_{nl2}	What is the uniform color of the referee?	Easy
Q_{nl3}	Is there any referee in the image?	Easy
Q_{nl4}	Which team does the goalkeeper belong to?	Medium
Q_{nl5}	Who is the defending team?	Medium
Q_{nl6}	Which part of the field are the players being now?	Hard
Q_{nl7}	How many players are there in the image?	Hard
Q_{nl8}	Is this image about corner kick? (If not, just list the correct one.)	??
Q_{nl9}	Is this image about penalty kick? (If not, just list the correct one.)	??
Q_{nl10}	Is this image about kick off? (If not, just list the correct one.)	??

Table 2: A set of questions

6.1. Performance of Visual Task Selecting Policy

To test the validity of reinforcement learning of selecting visual task modules, we test the inference time over accuracy with the state-of-art [2] [14] which is shown in Figure ??.

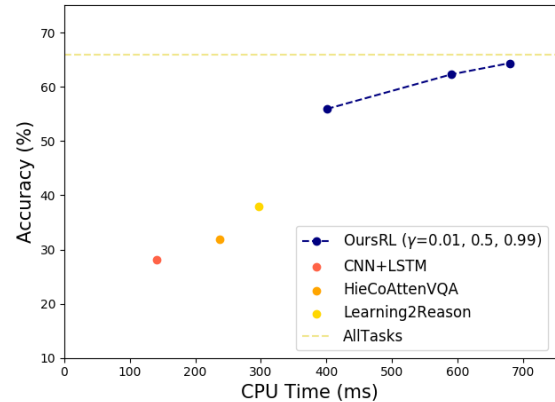


Figure 6: Inference Time and Accuracy

Here to test the generalization, we enlarge the training set by more various question with same meaning. For instance, the original question of Q_{nl5} is "Who is the defending team?", we add three more similar question asking "Who is attacking team?", "What is the uniform color of the defending team?" and "What is the uniform color of

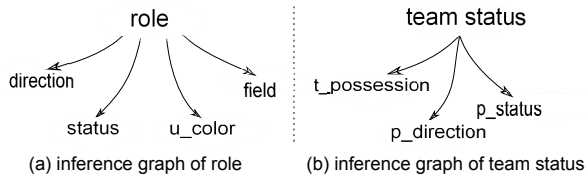


Figure 7: Inference Graphs

the attacking team?”. Unlike state-of-art methods answering questions in [25], adding generalization and variation in question would not dramatically change the performance, it is because the structure is not fixed, all the visual task selection is query oriented. For [10], even though the network is not fixed, the answering part is based on neural network, and essentially it also learns the statically correlation, which leads to the weakness in logical reasoning.

6.2. Effectiveness of Inference

Accuracy of Role. Only report results with Reinforcement learning

Accuracy of Team-Status. Only report results with Reinforcement learning

Accuracy of Kick-Off. [SHOW INFERENCE GRAPH AND RESULT TABLE!](#)

Accuracy of Penalty Kick. [SHOW INFERENCE GRAPH AND RESULT TABLE!](#)

Accuracy of Corner Kick. [SHOW INFERENCE GRAPH AND RESULT TABLE!](#)

Accuracy of Attacking Free Kick. [SHOW INFERENCE GRAPH AND RESULT TABLE!](#)

Accuracy of Balls. [Over new dataset. SHOW INFERENCE GRAPH AND RESULT TABLE!](#)

6.3. Overall Performance

We use the same question setting as [25], compared the following state-of-the-art methods: [2], and [10] with our method ($\gamma=0.99$), the overall performance is shown in Table 3.

References

- [1] J. Andreas, M. Rohrbach, T. Darrell, and D. Klein. Neural module networks. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun 2016. 2
- [2] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. L. Zitnick, and D. Parikh. VQA: Visual Question Answering. In *International Conference on Computer Vision (ICCV)*, 2015. 6, 7

	CNN+LSTM	HieCoAtten	Learn2Reason	Ours
Q_{nl1}	44.23	43.62	31.12	64.16
Q_{nl2}	71.31	77.66	9.4	47.43
Q_{nl3}	74.58	83.78	83.21	70.02
Q_{nl4}	40.48	39.29	51.92	62.14
Q_{nl5}	49.19	49.90	30.78	93.33
Q_{nl6}	20.56	18.70	30.0	62.13
Q_{nl7}	11.08	12.63	36.69	47.45
Avg.	46.40	49.11	51.08	65.97

Table 3: Accuracy comparison per query (%)

- [3] D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align and translate. May 2016. 3
- [4] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *CVPR*, 2017. 4
- [5] J. Dai, Y. Li, K. He, and J. Sun. R-FCN: object detection via region-based fully convolutional networks. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pages 379–387, 2016. 4
- [6] N. Friedman, D. Geiger, and M. Goldszmidt. Bayesian network classifiers. *Machine learning*, 29(2-3):131–163, 1997. 5
- [7] A. Fukui, D. H. Park, D. Yang, A. Rohrbach, T. Darrell, and M. Rohrbach. Multimodal compact bilinear pooling for visual question answering and visual grounding. *arXiv preprint arXiv:1606.01847*, 2016. 2
- [8] H. Gao, J. Mao, J. Zhou, Z. Huang, L. Wang, and W. Xu. Are you talking to a machine? dataset and methods for multilingual image question. In *Advances in neural information processing systems*, pages 2296–2304, 2015. 1
- [9] B. Goodrich and I. Arel. Reinforcement learning based visual attention with application to face detection. In *2012 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops, Providence, RI, USA, June 16-21, 2012*, pages 19–24, 2012. 2
- [10] R. Hu, J. Andreas, M. Rohrbach, T. Darrell, and K. Saenko. Learning to reason: End-to-end module networks for visual question answering. *CoRR*, abs/1704.05526, 3, 2017. 2, 7
- [11] D. Koller, N. Friedman, and F. Bach. *Probabilistic graphical models: principles and techniques*. MIT press, 2009. 5
- [12] S. Liu, Z. Zhu, N. Ye, S. Guadarrama, and K. Murphy. Improved image captioning via policy gradient optimization of spider. In *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017. 5
- [13] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. E. Reed, C. Fu, and A. C. Berg. SSD: single shot multibox detector. In *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part I*, pages 21–37, 2016. 4
- [14] J. Lu, J. Yang, D. Batra, and D. Parikh. Hierarchical question-image co-attention for visual question answering, 2016. 6

- [15] M. Malinowski, M. Rohrbach, and M. Fritz. Ask your neurons: A neural-based approach to answering questions about images. In *Proceedings of the IEEE international conference on computer vision*, pages 1–9, 2015. 1
- [16] S. Mathe, A. Pirinen, and C. Sminchisescu. Reinforcement learning for visual object detection. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 2894–2902, 2016. 2
- [17] V. Mnih, N. Heess, A. Graves, and K. Kavukcuoglu. Recurrent models of visual attention. In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, pages 2204–2212, 2014. 2
- [18] F. Petitjean, W. Buntine, G. I. Webb, and N. Zaidi. Accurate parameter estimation for bayesian network classifiers using hierarchical dirichlet processes. *Machine Learning*, 107(8-10):1303–1331, 2018. 5
- [19] M. Ren, R. Kiros, and R. Zemel. Image question answering: A visual semantic embedding model and a new dataset. *Proc. Advances in Neural Inf. Process. Syst.*, 1(2):5, 2015. 1
- [20] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 1, NIPS’15*, pages 91–99, Cambridge, MA, USA, 2015. MIT Press. 4
- [21] Z. Ren, X. Wang, N. Zhang, X. Lv, and L. Li. Deep reinforcement learning-based image captioning with embedding reward. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 1151–1159, 2017. 2
- [22] D. Teney, L. Liu, and A. van den Hengel. Graph-structured representations for visual question answering. *arXiv preprint*, 2017. 2
- [23] K. Tu, M. Meng, M. W. Lee, T. E. Choe, and S. C. Zhu. Joint video and text parsing for understanding events and answering queries. *IEEE MultiMedia*, 21(2):42–70, 2014. 2
- [24] S.-E. Wei, V. Ramakrishna, T. Kanade, and Y. Sheikh. Convolutional pose machines. In *CVPR*, 2016. 4
- [25] P. Xiong, H. Zhan, X. Wang, B. Sinha, and Y. Wu. Visual query answering by entity-attribute graph matching and reasoning. In *CVPR*, 2019. 2, 4, 5, 7
- [26] H. Xu and K. Saenko. Ask, attend and answer: Exploring question-guided spatial attention for visual question answering. In *European Conference on Computer Vision*, pages 451–466. Springer, 2016. 2
- [27] C. Zhu, Y. Zhao, S. Huang, K. Tu, and Y. Ma. Structured attentions for visual question answering. In *Proc. IEEE Int. Conf. Comp. Vis.*, volume 3, 2017. 2
- [28] Y. Zhu, O. Groth, M. Bernstein, and L. Fei-Fei. Visual7w: Grounded question answering in images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4995–5004, 2016. 2