

A Focused Dynamic Attention Model for Visual Question Answering

Ilija Ilievski, Shuicheng Yan, Jiashi Feng

ilija.ilievski@u.nus.edu, {eleyans, elefjia}@nus.edu.sg

National University of Singapore

Abstract. Visual Question and Answering (VQA) problems are attracting increasing interest from multiple research disciplines. Solving VQA problems requires techniques from both computer vision for understanding the visual contents of a presented image or video, as well as the ones from natural language processing for understanding semantics of the question and generating the answers. Regarding visual content modeling, most of existing VQA methods adopt the strategy of extracting global features from the image or video, which inevitably fails in capturing fine-grained information such as spatial configuration of multiple objects. Extracting features from auto-generated regions – as some region-based image recognition methods do – cannot essentially address this problem and may introduce some overwhelming irrelevant features with the question. In this work, we propose a novel Focused Dynamic Attention (FDA) model to provide better aligned image content representation with proposed questions. Being aware of the key words in the question, FDA employs off-the-shelf object detector to identify important regions and fuse the information from the regions and global features via an LSTM unit. Such question-driven representations are then combined with question representation and fed into a reasoning unit for generating the answers. Extensive evaluation on a large-scale benchmark dataset, VQA, clearly demonstrate the superior performance of FDA over well-established baselines.

Keywords: Visual Question Answering, Attention

1 Introduction

Visual question answering (VQA) is an active research direction that lies in the intersection of computer vision, natural language processing, and machine learning. Even though with a very short history, it already has received great research attention from multiple communities. Generally, the VQA investigates a generalization of traditional QA problems where visual input (*e.g.*, an image) is necessary to be considered. More concretely, VQA is about how to provide a correct answer to a human posed question concerning contents of one presented image or video.

VQA is a quite challenging task and undoubtedly important for developing modern AI systems. The VQA problem can be regarded as a Visual Turing Test

[1,2], and besides contributing to the advancement of the involved research areas, it has other important applications, such as blind person assistance and image retrieval. Coming up with solutions to this task requires natural language processing techniques for understanding the questions and generating the answers, as well as computer vision techniques for understanding contents of the concerned image. With help of these two core techniques, the computer can perform reasoning about the perceived contents and posed questions.

Recently, VQA is advanced significantly by the development of machine learning methods (in particular the deep learning ones) that can learn proper representations of questions and images, align and fuse them in a joint question-image space and provide a direct mapping from this joint representation to a correct answer.

For example, consider the following image-question pair: an image of an apple tree with a basket of apples next to it, and a question “How many apples are in the basket?”. Answering this question requires VQA methods to first understand the semantics of the question, then locate the objects (apples) in the image, understand the relation between the image objects (which apples are in the basket), and finally count them and generate an answer with the correct number of apples.

The first feasible solution to VQA problems was provided by Malinowski and Fritz in [2], where they used a semantic language parser and a Bayesian reasoning model, to understand the meaning of questions and to generate the proper answers. Malinowski and Fritz also constructed the first VQA benchmark dataset, named as DAQUAR, which contains 1,449 images and 12,468 questions generated by humans or automatically by following a template and extracting facts from a database [2]. Shortly after, Ren et al. [3] released the TORONTO-QA dataset, which contains a large number of images (123,287) and questions (117,684), but the questions are automatically generated and thus can be answered without complex reasoning. Nevertheless, the release of the TORONTO-QA dataset was important since it provided enough data for deep learning models to be trained and evaluated on the VQA problem [4,5,3]. More recently, Antol et al. [6] published the currently largest VQA dataset. It consists of three human posed questions and ten answers given by different human subjects, for each one of the 204,721 images found in the Microsoft COCO dataset [7]. Answering the 614,163 questions requires complex reasoning, common sense, and real-world knowledge, making the VQA dataset suitable for a true Visual Turing Test. The VQA authors split the evaluation on their dataset on two tasks: an open-ended task, where the method should generate a natural language answer, and a multiple-choice task, where for each question the method should chose one of the 18 different answers.

The current top performing methods [8,9,10] employ deep neural network model that predominantly uses the convolutional neural network (CNN) architecture [11,12,13,14] to extract image features and a Long Short-Term Memory (LSTM) [15] network to extract the representations for questions. The CNN and LSTM representation vectors are then usually fused by concatenation [16,3,5] or

element-wise multiplication [17,18]. Other approaches additionally incorporate some kind of attention mechanism over the image features [18,19,20].

Properly modeling the image contents is one of the critical factors for solving VQA problems well. A common practice with existing VQA methods on modeling image contents is to extract global features for the overall image. However, only using global feature is arguably insufficient to capture all the necessary visual information and provide full understanding of image contents such as multiple objects, spatial configuration of the objects and informative background. This issue can be relieved to some extent by extracting features from object proposals – the image regions that possibly contain objects of interest. However, using features from all image regions [19,18] may provide too much noise or overwhelming information irrelevant to the question and thus hurt the overall VQA performance.

In this work, we propose a question driven attention model that is able to automatically identify and focus on image regions relevant for the current question. We name our proposed model Focused Dynamic Attention (FDA) for Visual Question Answering. With the FDA model, computers can select and recognize the image regions in a well-aligned sequence with the key words containing in a given question. Recall the above VQA example. To answer the question of “How many apples are in the basket?”, FDA would first localize the regions corresponding to the key words “apples” and “basket” (with the help of a generic object detector) and extract description features from these regions of interest. Then VQA compliments the features from selected image regions with a global image feature providing contextual information for the overall image, and reconstruct a visual representation by encoding them with a Long Short-Term Memory (LSTM) unit.

We evaluate and compare the performance of our proposed FDA model on two types of VQA tasks, *i.e.*, the open-ended task and the multiple-choice task, on the VQA dataset – the largest VQA benchmark dataset. Extensive experiments demonstrate that FDA brings substantial performance improvement upon well-established baselines.

The main contributions of this work can be summarized as follows:

- We introduce a focused dynamic attention mechanism that learns to use the question word order to shift the focus from one image object, to another.
- We describe a model that fuses local and global context visual features with textual features.
- We perform an extensive evaluation, comparing to all existing methods, and achieve state-of-the-art accuracy on the open-ended, and on the multiple-choice VQA tasks.

The rest of the paper is organized as follows. In Section 2 we review the current VQA models, and compare them to our model. We formulate the problem and explain our motivation in Section 3. We describe our model in Section 4 and in Section 5 we evaluate and compare it with the current state-of-the-art models. We conclude our work in Section 6.

2 Related Work

VQA has received great research attention recently and a couple of methods have been developed to solve this problem. The most similar model to ours is the Stacked Attention Networks (SAN) proposed by Yang et al. [19]. Both models use attention mechanism that combines the words and image regions. However, [19] use convolutional neural network to put attention over the image regions, based on the question word unigrams, bigrams, and trigrams. Further, their attention mechanism is not using object bounding boxes, which makes the attention less focused.

Another model that uses attention mechanism in solving VQA problems is the ABC-CNN model described in [18]. ABC-CNN uses the question embedding to configure convolutional kernels that will define an attention weighted map over the image features. The advantage of our FDA model over ABC-CNN is two fold. First, FDA employs an LSTM network to encode the image region features in a order that corresponds to the question word order. Second, FDA does not put handcrafted weights on the image features (a practice showed to hurt the learning process in our experiments). Instead, FDA extracts CNN features directly from the cropped image regions of interest. In this sense, FDA is more efficient than ABC-CNN in visual contents modeling.

Yet another attention model for visual question answering is proposed in [17]. The work, is closely related to the work by [18], in that it also applies a weighted map over the image and the question word features. However, similar to our work, they use object proposals from [21] to select image regions instead of the whole image. Different from that work, our proposed FDA model also employs the information embedded in the order of the question words and focuses on the corresponding object bounding boxes. In contrast, the model proposed in [21] straightforwardly concatenate all the image region features with the question word features and feed them all at once to a two layer network.

Jiang et al. propose another model that combines the CNN image features and an LSTM network for encoding the multimodal representation, with the addition of a Compositional Memory units which fuse the image and word feature vectors [22].

Ma et al. in [10] take an interesting approach and use three convolutional neural networks to represent not only the image, but also the question, and their common representation in a multimodal space. The multimodal representation is then fed to a SoftMax layer to produce the answer.

Another interesting approach worth mentioning is the work by Andreas et al. [23]. They use a semantic grammar parser to parse the question and propose neural network layouts accordingly. They train a model to learn to compose a network from one of the proposed network layouts using several types of neural modules, each specifically designed to address the different sub-tasks of the VQA problem (e.g. counting, locating an object, etc.).

3 Method Overview

In this section, we briefly describe the motivation and give formal problem formulation.

3.1 Problem Formulation

The visual question answering problem can be represented as predicting the best answer \hat{a} given an image I and a question q . Common practice [6,16,3,19] is to use the 1,000 most common answers in the training set and thus simplify the VQA task to a classification problem. The following equation represents the problem mathematically:

$$\hat{a} = \arg \max_{a \in \Omega} p(a|I, q; \theta) \quad (1)$$

where Ω is the set of all possible answers and θ are the model weights.

3.2 Motivation

The baseline methods from [6] show only modest increase in accuracy when including the image features (4.98% for open-ended questions, and 2.42% for multiple-choice question). We believe that the image contains a lot more information and should increase the accuracy much more. Thus, we focus on improving the image features and design a visual attention mechanism, which learns to focus on the question related image regions.

The proposed attention mechanism is loosely inspired on the human visual attention mechanism. Humans shift the focus from one image region to another, before understanding how the regions relate to each other and grasping the meaning of the whole image. Similarly, we feed our model image regions relevant for the question at hand, before showing the whole image.

4 Focused Dynamic Attention for VQA

The FDA model is composed of question and image understanding components, attention mechanism, and a multimodal representation fusion network (Figure 1). In this section we describe them individually.

4.1 Question Understanding

Following a common practice, our FDA model uses an LSTM network to encode the question in a vector representation [15,5,18,8]. The LSTM network learns to keep in its state the feature vectors of the important question words, and thus provides the question understanding component with a word attention mechanism.

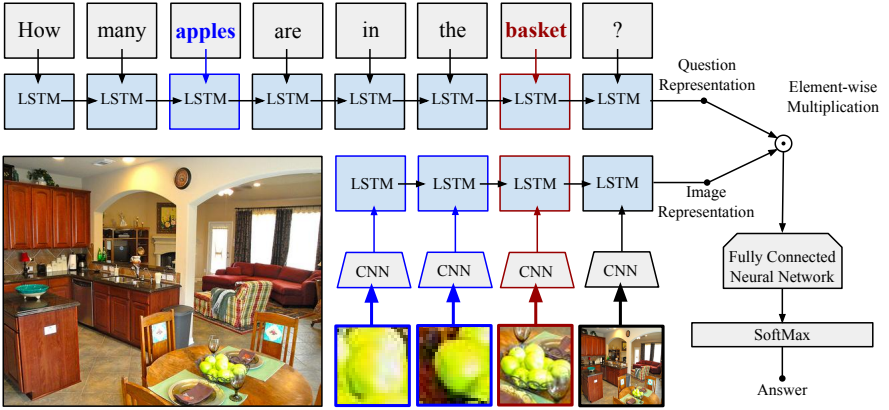


Fig. 1. Focused dynamic attention model diagram.

4.2 Image Understanding

Following prior work [3,4,5], we use a pre-trained convolutional neural network (CNN) to extract image feature vectors. Specifically, we use the Deep Residual Networks model used in ILSVRC and COCO 2015 competitions, which won the 1st places in: ImageNet classification, ImageNet detection, ImageNet localization, COCO detection, and COCO segmentation [24]. We extract the weights of the layer immediately before the final SoftMax layer and regard them as visual features. We extract such features for the whole image (global visual features) and for the specific image regions (local visual features). However, contrary to the existing approaches, we employ an LSTM network to combine the local and global visual features into a joint representation.

4.3 Focused Dynamic Attention Mechanism

We introduce a focused dynamic attention mechanism that learns to focus on image regions related to the question words.

The attention mechanism works as follows. For each image object¹ it uses word2vec word embeddings [26] to measure the similarity between the question words and the object label. Next, it selects objects with similarity score greater than 0.5 and extracts the feature vectors of the objects bounding boxes with a pre-trained ResNet model [24]. Following the question word order, it feeds the LSTM network with the corresponding object feature vectors. Finally, it feeds the LSTM network with the feature vector of the whole image and it uses the resulting LSTM state as a visual representation. Thus, the attention mechanism

¹ During training we use the ground truth object bounding boxes and labels. At test time we use the precomputed bounding boxes from [25] and classify them with [24] to obtain the object labels.

enables the model to combine the local and global visual features into a single representation, necessary for answering complex visual questions.

Figure 1 illustrates the focused dynamic attention mechanism with an example.

4.4 Multimodal Representation Fusion

We regard the final state of the two LSTM networks as a question and image representation. We start fusing them into single representation by applying *Tanh* on the question representation and *ReLU*² on the image representation³. We proceed by doing an element-wise multiplication of the two vector representations and the resulting vector is fed to a fully-connected neural network. Finally a SoftMax layer classify the multimodal representation into one of the possible⁴ answers.

5 Evaluation

In this section we detail the model implementation and compare our model against the current state-of-the-art methods.



					
Why does this male have his arms in this position?	balance for balance for balance	angry he's carrying bags hug	How many people are wearing an orange shirt?	3 3 3	1 3 3
Are the clouds high in the sky?	yes yes yes	no no yes	Is this a trained elephant?	yes yes yes	yes yes yes

Fig. 2. Representative examples of questions (black), (a subset of the) answers given when looking at the image (green), and answers given when not looking at the image (blue) for two images from the VQA dataset. Examples provided by [6].

² Defined as $f(x) = \max(0, x)$.

³ Applying different activation functions gave slightly worse overall results

⁴ We follow [6] and use the 1000 most common answers

5.1 Dataset

For all experiments we use the Visual Question Answering (VQA) dataset [6], which is the largest and most complex image dataset for the visual question answering task. The dataset contains three human posed questions and ten answers given by different human subjects, for each one of the 204,721 images found in the Microsoft COCO dataset [7]. Figure 2 shows two representative examples found in the dataset. The evaluation is done on following two test splits test-dev and test-std and on following two tasks:

- An open-ended task, where the method should generate a natural language answer;
- A multiple-choice task, where for each question the method should chose one of the 18 different answers.

We evaluate the performance of all the methods in the experiments using the public evaluation server for fair evaluation.

5.2 Baseline Model

We compare our model against the baseline models provided by the VQA dataset authors [27], which currently achieve the best performance on the test-standard split for the multiple-choice task. The model, first described in [6], is a standard implementation of an LSTM+CNN VQA model. It uses an LSTM to encode the question and CNN features to encode the image. To answer a question it multiplies the last LSTM state with the image CNN features and feeds the result into a SoftMax layer for classification into one of the 1,000 most common answers. The implementation in [27] uses a deeper two layer LSTM network for encoding the question, and normalized image CNN features, which showed crucial for achieving the state-of-the-art.

5.3 Model Implementation and Training Details

We transform the question words into a vector form by multiplying one-hot vector representation with a word embedding matrix. The vocabulary size is 12,602 and the word embeddings are 300 dimensional. We feed a pre-trained ResNet network [24] and use the 2,048 dimensional weight vector of the layer before the last fully-connected layer.

The word and image vectors are feed into two separate LSTM networks. The LSTM networks are standard implementation of one layer LSTM network [15], with a 512 dimensional state vector. The final state of the question LSTM is passed through $Tanh$, while the final state of the image LSTM is passed through $ReLU$ ⁵. We do element-wise multiplication on the resulting vectors, to obtain a multimodal representation vector, which is then fed to a fully-connected neural network.

⁵ Defined as $f(x) = \max(0, x)$.

Table 1. Comparison between the baselines from [6], the state-of-the-art models and our FDA model on VQA test-dev and test-standard data for the **open-ended** task. Results from most recent methods including CM [22], ACK [9], iBOWIMG [16], DPPnet [8], D-NMN [23], D-LSTM [27], and SAN [19] are provided and compared with.

Method	test-dev				test-std
	All	Y/N	Other	Num	All
VQA					
Question	48.09	75.66	27.14	36.70	-
Image	28.13	64.01	3.77	0.42	-
Q+I	52.64	75.55	37.37	33.67	-
LSTM Q+I	53.74	78.94	36.42	35.24	54.06
CM	52.62	78.33	35.93	34.46	
ACK	55.72	79.23	40.08	36.13	55.98
iBOWIMG	55.72	76.55	42.62	35.03	55.89
DPPnet	57.22	80.71	41.71	37.24	57.36
D-NMN	57.90	80.50	43.10	37.40	58.00
D-LSTM	-	-	-	-	58.16
SAN	58.70	79.30	46.10	36.6	58.90
FDA	59.24	81.14	45.77	36.16	59.54

Table 2. Comparison between the baselines from [6], the state-of-the-art models and our FDA model on VQA test-dev and test-standard data for the **multiple-choice** task. Results from most recent methods including WR [17], iBOWIMG [16], DPPnet [8], and D-LSTM [27] are also shown for comparison.

Method	test-dev				test-std
	All	Y/N	Other	Num	All
VQA					
Question	53.68	75.71	38.64	37.05	-
Image	30.53	69.87	3.76	0.45	-
Q+I	58.97	75.59	50.33	34.35	-
LSTM Q+I	57.17	78.95	43.41	35.80	57.57
WR	60.96	-	-	-	-
iBOWIMG	61.68	76.68	54.44	38.94	61.97
DPPnet	62.48	80.79	52.16	38.94	62.69
D-LSTM	-	-	-	-	63.09
FDA	64.01	81.50	54.72	39.00	64.18

5.4 Model Evaluation and Comparison

We compare our model with the baselines provided by the VQA authors [6]. The results for the open-ended task are listed in Table 1 and the results for the multiple-choice task are given in Table 2. In the tables, the “Question” and “Image” baselines are only using the question words and the image, respectively. The “Q+I” is a baseline that combines the two, but do not use an LSTM network. “LSTM Q+I” and “D-LSTM” are LSTM models, with one and two layers ac-

cordingly. Comparing the performance of baselines we can observe the accuracy increase with the addition of information from each modality.

From Table 1, one can observe that our proposed FDA model achieves the best performance on this benchmark dataset. It outperforms the state-of-the-art (SAN) with a margin of around 0.6%. The SAN model also employs attention to focus on specific regions. However, their attention model (without access to the automatically generated bounding boxes) is focusing on more spread regions which may include cluttered and noisy background. In contrast, FDA only focuses on the selected regions and extracts more clean information for answering the questions. This is the main reason that FDA can outperform SAN although these two methods are both based on attention models.

The advantage of employing focused dynamic attention in FDA is more significant when solving the multiple-choice VQA problems. From Table 2, one can observe that our proposed FDA model achieves the best ever performance on the VQA dataset. In particular, it improves the performance of the state-of-the-art (D-LSTM) by a margin of 1.1% which is quite significant for this challenging task. The D-LSTM method employs a deeper network to enhance the discriminative capacity of the visual features. However, they do not identify the informative regions for answering the questions. In contrast, FDA incorporates the automatic region localization by employing a question-driven attention model. This is helpful for filtering out irrelevant noise, and establishing the correspondence between regions and candidate answers. Thus FDA gives substantial performance improvement.

5.5 Qualitative Results

We qualitatively evaluate our model on a set of examples where complex reasoning and focusing on the relevant local visual features are needed for answering the question correctly.

Figure 3 shows particularly difficult examples (the predominant image color is **not** the correct answer) of “What color” type of questions. But, by focusing on the question related image regions, the FDA model is still able to produce the correct answer.

In Figure 4 we show examples where the model focuses on different regions from the **same** image, depending on the words in the question. Focusing on the right image region is crucial when answering unusual questions for an image (Row 1), questions about small image objects (Row 2), or when the most dominant image object partly occludes the question related region and can lead to a wrong answer (Row 3).

Representative examples of questions that require image object identification are shown in Figure 5. We can observe that the focused attention enables the model to answer complex questions (Row 1, left) and counting questions (Row 1, right). The question guided image object identification greatly simplifies the answering of questions like the ones shown in Row 2 and Row 3.



What color is the **frisbee**?
- Red.



What color are the **glass** items?
- Green.



What color is the **mouse**?
- White.



What color is the lady's **umbrella**?
- Blue

Fig. 3. Representative examples where focusing on the question related objects helps FDA answer “What color” type of questions. The question words in bold have been matched with an image region. The yellow region caption box contains the question word, followed by the region label, and in parenthesis their cosine similarity (see Section 4.3 for more details).



Is this a birthday **cake**?
- Yes.



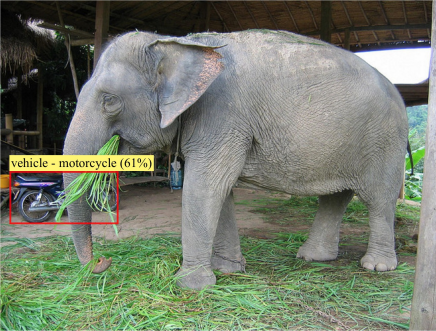
Is someone in all likelihood, a **zoo** fancier? - Yes.



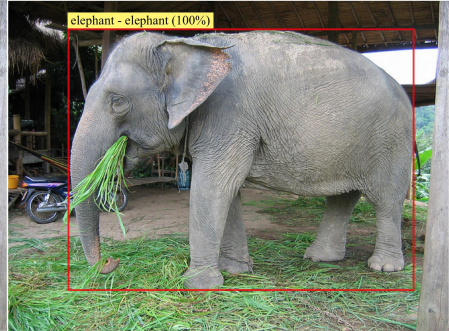
What **fruit** is by the **sink**?
- Apples.



Is there a **cookbook** in the picture?
- Yes.

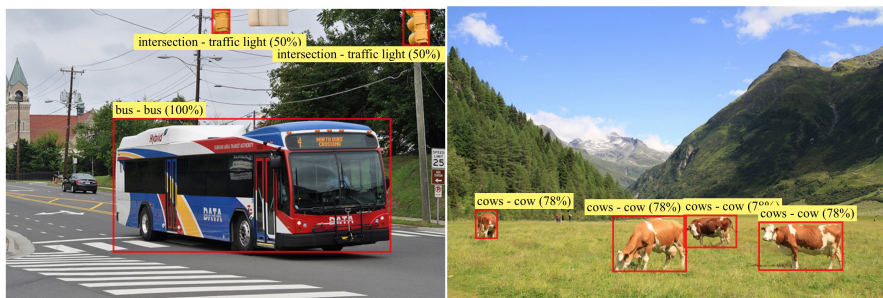


What type of **vehicle** is pictured?
- Motorcycle.



Does the **elephant** have tusks?
- No.

Fig. 4. Representative examples where the model focuses on different regions from the **same** image, depending on the question. The question words in **bold** have been matched with an image region. The yellow region caption box contains the question word, followed by the region label, and in parenthesis their cosine similarity (see Section 4.3 for more details).



Is the **bus** in the middle of the **intersection**? - Yes.

How many **cows** are present?
- 4.



Does this **dessert** have any **fruit** with it? - Yes.

Is this a **carrot cake**?
- Yes.



Is this an **airport**?
- Yes.

Is the individual skiing or **snowboarding**? - Snowboarding.

Fig. 5. Representative examples of questions that require image object identification. The question words in bold have been matched with an image region. The yellow region caption box contains the question word, followed by the region label, and in parenthesis their cosine similarity (see Section 4.3 for more details).

6 Conclusion

In this work, we proposed a novel Focused Dynamic Attention (FDA) model to solve the challenging VQA problems. FDA is built upon a generic object-centric attention model for extracting question related visual features from an image as well as a stack of multiple LSTM layers for feature fusion. By only focusing on the identified regions specific for proposed questions, FDA was shown to be able to filter out overwhelming irrelevant informations from cluttered background or other regions, and thus substantially improved the quality of visual representations in the sense of answering proposed questions. By fusing cleaned regional representation, global context and question representation via LSTM layers, FDA provided significant performance improvement over baselines on the VQA benchmark datasets, for both the open-ended and multiple-choices VQA tasks. Excellent performance of FDA clearly demonstrates its stronger ability of modeling visual contents and also verifies paying more attention to visual part in VQA tasks could essentially improve the overall performance. In the future, we are going to further explore along this research line and investigate different attention methods for visual information selection as well as better reasoning model for interpreting the relation between visual contents and questions.

References

1. Geman, D., Geman, S., Hallonquist, N., Younes, L.: Visual turing test for computer vision systems. *Proceedings of the National Academy of Sciences* **112**(12) (2015) 3618–3623
2. Malinowski, M., Fritz, M.: A multi-world approach to question answering about real-world scenes based on uncertain input. In: *Advances in Neural Information Processing Systems*. (2014) 1682–1690
3. Ren, M., Kiros, R., Zemel, R.: Exploring models and data for image question answering. In: *Advances in Neural Information Processing Systems*. (2015) 2935–2943
4. Gao, H., Mao, J., Zhou, J., Huang, Z., Wang, L., Xu, W.: Are you talking to a machine? dataset and methods for multilingual image question. In: *Advances in Neural Information Processing Systems*. (2015) 2287–2295
5. Malinowski, M., Rohrbach, M., Fritz, M.: Ask your neurons: A neural-based approach to answering questions about images. In: *Proceedings of the IEEE International Conference on Computer Vision*. (2015) 1–9
6. Antol, S., Agrawal, A., Lu, J., Mitchell, M., Batra, D., Lawrence Zitnick, C., Parikh, D.: Vqa: Visual question answering. In: *The IEEE International Conference on Computer Vision (ICCV)*. (December 2015)
7. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: *Computer Vision—ECCV 2014*. Springer (2014) 740–755
8. Noh, H., Seo, P.H., Han, B.: Image question answering using convolutional neural network with dynamic parameter prediction. *arXiv preprint arXiv:1511.05756* (2015)

9. Wu, Q., Wang, P., Shen, C., Hengel, A.v.d., Dick, A.: Ask me anything: Free-form visual question answering based on knowledge from external sources. arXiv preprint arXiv:1511.06973 (2015)
10. Ma, L., Lu, Z., Li, H.: Learning to answer questions from image using convolutional neural network. arXiv preprint arXiv:1506.00333 (2015)
11. LeCun, Y., Boser, B., Denker, J.S., Henderson, D., Howard, R.E., Hubbard, W., Jackel, L.D.: Backpropagation applied to handwritten zip code recognition. *Neural computation* **1**(4) (1989) 541–551
12. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: *Advances in neural information processing systems*. (2012) 1097–1105
13. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556 (2014)
14. Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., Rabinovich, A.: Going deeper with convolutions. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. (2015) 1–9
15. Hochreiter, S., Schmidhuber, J.: Long short-term memory. *Neural computation* **9**(8) (1997) 1735–1780
16. Zhou, B., Tian, Y., Sukhbaatar, S., Szlam, A., Fergus, R.: Simple baseline for visual question answering. arXiv preprint arXiv:1512.02167 (2015)
17. Shih, K.J., Singh, S., Hoiem, D.: Where to look: Focus regions for visual question answering. arXiv preprint arXiv:1511.07394 (2015)
18. Chen, K., Wang, J., Chen, L.C., Gao, H., Xu, W., Nevatia, R.: Abc-cnn: An attention based convolutional neural network for visual question answering. arXiv preprint arXiv:1511.05960 (2015)
19. Yang, Z., He, X., Gao, J., Deng, L., Smola, A.: Stacked attention networks for image question answering. arXiv preprint arXiv:1511.02274 (2015)
20. Xu, H., Saenko, K.: Ask, attend and answer: Exploring question-guided spatial attention for visual question answering. arXiv preprint arXiv:1511.05234 (2015)
21. Zitnick, C.L., Dollár, P.: Edge boxes: Locating object proposals from edges. In: *Computer Vision–ECCV 2014*. Springer (2014) 391–405
22. Jiang, A., Wang, F., Porikli, F., Li, Y.: Compositional memory for visual question answering. arXiv preprint arXiv:1511.05676 (2015)
23. Andreas, J., Rohrbach, M., Darrell, T., Klein, D.: Learning to compose neural networks for question answering. arXiv preprint arXiv:1601.01705 (2016)
24. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. arXiv preprint arXiv:1512.03385 (2015)
25. Pont-Tuset, J., Arbeláez, P., Barron, J., Marques, F., Malik, J.: Multiscale combinatorial grouping for image segmentation and object proposal generation. In: arXiv:1503.00848. (March 2015)
26. Mikolov, T., Sutskever, I., Chen, K., Corrado, G.S., Dean, J.: Distributed representations of words and phrases and their compositionality. In: *Advances in neural information processing systems*. (2013) 3111–3119
27. Jiasen Lu, Xiao Lin, D.B., Parikh, D.: Deeper lstm and normalized cnn visual question answering model. https://github.com/VT-vision-lab/VQA_LSTM_CNN (2015)