# Improved Fusion of Visual and Language Representations by Dense Symmetric Co-Attention for Visual Question Answering

Duy-Kien Nguyen<sup>1</sup> and Takayuki Okatani<sup>1,2</sup>
<sup>1</sup>Tohoku University <sup>2</sup>RIKEN Center for AIP

{kien, okatani}@vision.is.tohoku.ac.jp

# **Abstract**

A key solution to visual question answering (VQA) exists in how to fuse visual and language features extracted from an input image and question. We show that an attention mechanism that enables dense, bi-directional interactions between the two modalities contributes to boost accuracy of prediction of answers. Specifically, we present a simple architecture that is fully symmetric between visual and language representations, in which each question word attends on image regions and each image region attends on question words. It can be stacked to form a hierarchy for multi-step interactions between an image-question pair. We show through experiments that the proposed architecture achieves a new state-of-the-art on VQA and VQA 2.0 despite its small size. We also present qualitative evaluation, demonstrating how the proposed attention mechanism can generate reasonable attention maps on images and questions, which leads to the correct answer prediction.

### 1. Introduction

There has been a significant progress in the study of visual question answering (VQA) over a short period of time since its introduction, showing rapid boost of performance for common benchmark datasets. This progress has been mainly brought about by two lines of research, the development of better attention mechanisms and the improvement in fusion of features extracted from an input image and question.

Since introduced by Bahdanau et al. [3], attention has been playing an important role in solutions of various problems of artificial intelligence ranging from tasks using single modality (e.g., language, speech, and vision) to multimodal tasks. For VQA, attention on image regions generated from the input question was first introduced [32] and then several extensions have been proposed [21, 35, 5]. Meanwhile, researchers have proposed several methods for

feature fusion [6, 16, 36], where the aim is to obtain better fused representation of image and question pairs. These studies updated the state-of-the-art for common benchmark datasets at the time of each publication.

We observe that these two lines of research have been independently conducted so far. This is particularly the case with the studies of feature fusion methods, where attention is considered to be optional, even though the best performance is achieved with it. However, we think that they are rather *two different approaches towards the same goal*. In particular, we argue that a better attention mechanism leads to a better fused representation of image-question pairs.

Motivated by this, we propose a novel co-attention mechanism for improved fusion of visual and language representations. Given representations of an image and a question, it first generates an attention map on image regions for each question word and an attention map on question words for each image region. It then performs computation of attended features, concatenation of multimodal representations, and their transformation by a single layer network with ReLU and a residual connection. These computations are encapsulated into a composite network that we call dense co-attention layer, since it considers every interaction between any image region and any question word. The layer has fully symmetric architecture between the two modalities, and can be stacked to form a hierarchy that enables multi-step interactions between the image-question pair.

Starting from initial representations of an input image and question, each dense co-attention layer in the layer stack updates the representations, which are inputted to the next layer. Its final output are then fed to a layer for answer prediction. We use additional attention mechanisms in the initial feature extraction as well as the answer prediction layer. We call the entire network including all these components the *dense coattention network* (DCN). We show the effectiveness of DCNs by several experimental results; they achieve the new state-of-the-art for VQA 1.0 and 2.0 datasets.

## 2. Related Work

In this section, we briefly review previous studies of VQA with a special focus on the developments of attention mechanisms and fusion methods.

### 2.1. Attention Mechanisms

Attention has proved its effectiveness on many tasks and VQA is no exception. A number of methods have been developed so far, in which question-guided attention on image regions is commonly used. They are categorized into two classes according to the type of employed image features. One is the class of methods that use visual features from some region proposals, which are generated by Edge Boxes [26, 12] or Region Proposal Network [28]. The other is the class of methods that use convolutional features (i.e., activations of convolutional layers) [5, 6, 14, 15, 16, 21, 23, 24, 31, 32, 36].

There are several approaches to creation and use of attention maps. Yang *et al.* [32] developed stacked attention network that produces multiple attention maps on the image in a sequential manner, aiming at performing multiple steps of reasoning. Kim *et al.* [15] extended this idea by incorporating it into a residual architecture to produce better attention information. Chen *et al.* [5] proposed a structured attention model that can encode cross-region relation, aiming at properly answering questions that involve complex inter-region relations.

Earlier studies mainly considered question-guided attention on image regions. In later studies, the opposite orientation of attention, i.e., image-guided attention on question words, is considered additionally. Lu *et al.* [21] introduced the co-attention mechanism that generates and uses attention on image regions and on question words. To reduce the gap of image and question features, Yu *et al.* [35] utilized attention to extract not only spatial information but also language concept of the image. Yu *et al.* [36] combined the mechanism with a novel multi-modal feature fusion of image and question.

We point out that the existing attention mechanisms only consider a limited amount of possible interactions between image regions and question words. Some consider only attention on image regions from a *whole* question. Coattention additionally considers attention on question words but it is created from a *whole* image. We argue that this can be a significant limitation of the existing approaches. The proposed mechanism can deal with every interaction between any image region and any question word, which possibly enables to model unknown complex image-question relations that are necessary for correctly answering questions.

#### 2.2. Multimodal Feature Fusion

The common framework of existing methods is that visual and language features are independently extracted from the image and question at the initial step, and they are fused at a later step to compute the final prediction. In early studies, researchers employed simple fusion methods such as the concatenation, summation, and element-wise product of the visual and language features, which are fed to fully connected layers to predict answers.

It was then shown by Fukui et al. [6] that a more complicated fusion method does improve prediction accuracy; they introduced the bilinear (pooling) method that uses an outer product of two vectors of visual and language features for their fusion. As the outer product gives a very high-dimensional feature, they adopt the idea of Gao et al. [7] to compress the fused feature and name it the Multimodal Compact Bilinear (MCB) pooling method. However, the compacted feature of the MCB method still tends to be high-dimensional to guarantee robust performance, Kim et al. [16] proposed low-rank bilinear pooling using Hadamard product of two feature vectors, which is called the Multimodal Low-rank Bilinear (MLB) pooling. Pointing out that MLB suffers from slow convergence rate, Yu et al. [36] proposed the Multi-modal Factorized Bilinear (MFB) pooling, which computes a fused feature with a matrix factorization trick to reduce the number of parameters and improve convergence rate.

The attention mechanisms can also be considered feature fusion methods, regardless of whether it is explicitly mentioned, since they are designed to obtain a better representation of image-question pairs based on their interactions. This is particularly the case with co-attention mechanisms in which the two features are treated symmetrically. Our dense co-attention network is based on this observation. It fuses the two features by multiple applications of the attention mechanism that can use more fine-grained interactions between them.

# 3. Dense Co-Attention Network (DCN)

In this section, we describe the architecture of DCNs; see Fig.1 for its overview. It consists of a stack of *dense co-attention layers* that fuses language and visual features repeatedly, on top of which an *answer prediction layer* that predict answers in a multi-label classification setting [28]. We first explain the initial feature extraction from the input question and image (Sec.3.1) and then describe the dense co-attention layer (Sec.3.2) and the answer prediction layer (Sec.3.3).

### 3.1. Feature Extraction

We employ pretrained networks that are commonly used in previous studies [15, 33, 16, 5] for encoding or extracting features from images, questions, and answers, such as

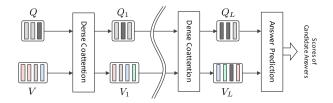


Figure 1: The global structure of the dense co-attention network (DCN).

pretrained ResNet [10] with some differences from earlier studies.

#### 3.1.1 Question and Answer Representation

We use bi-directional LSTM for encoding questions and answers. Specifically, a question consisting of N words is first converted into a sequence  $\{e_1^Q,...,e_N^Q\}$  of GloVe vectors [25], which are then inputted into a one-layer bi-directional LSTM (Bi-LSTM) with a residual connection as

$$\overrightarrow{q_n} = \text{Bi-LSTM}(\overrightarrow{q_{n-1}}, e_n^Q),$$
 (1)

$$\overleftarrow{q_n} = \text{Bi-LSTM}(\overleftarrow{q_{n+1}}, e_n^Q).$$
 (2)

We then create a matrix  $Q=[q_1,...,q_N]\in\mathbb{R}^{d\times N}$  where  $q_n=[\overrightarrow{q_n}^{\top},\overleftarrow{q_n}^{\top}]^{\top}$  (n=1,...,N). We will also use  $s_Q=[\overrightarrow{q_N}^{\top},\overleftarrow{q_1}^{\top}]^{\top}$ , concatenation of the last hidden states in the two paths, for obtaining representation of an input image (Sec.3.1.2). We randomly initialized the Bi-LSTM. It is worth noting that we initially used a pretrained two-layer Bi-LSTM that yields Context Vectors (CoVe) in [22], which we found does not contribute to performance.

We follow a similar procedure to encode answers. An answer of M words is converted into  $\{e_1^A,...,e_M^A\}$  and then inputted to the same Bi-LSTM, yielding the hidden states  $\overrightarrow{a_m}$  and  $\overleftarrow{a_m}$   $(m=1,\ldots,M)$ . We will use  $s_A=[\overrightarrow{a_M}^\top,\overleftarrow{a_1}^\top]^\top$  for answer representation.

### 3.1.2 Image Representation

As in many previous studies, we use a pretrained CNN (i.e., a ResNet [10] with 152 layers pretrained on ImageNet) to extract visual features of multiple image regions, but our extraction method is slightly different. We extract features from four conv. layers and then use a question-guided attention on these layers to fuse their features. We do this to exploit the maximum potential of the subsequent dense co-attention layers. We conjecture that features at different levels in the hierarchy of visual representation [37, 34] will be necessary to correctly answer a wide range of questions.

To be specific, we extract outputs from the four conv. layers (after ReLU) before the last four pooling layers. These are tensors of different sizes (i.e.,  $256 \times 112 \times 112$ ,  $512 \times 56 \times 56$ ,  $1024 \times 28 \times 28$ , and  $2048 \times 14 \times 14$ ) and

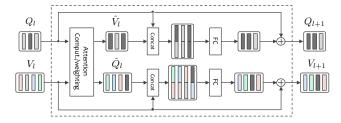


Figure 2: The internal structure of a single dense coattention layer of layer index l + 1.

are converted into tensors of the same size  $(d \times 14 \times 14)$  by applying max pooling with a different pooling size and one-by-one convolution to each. We also apply  $l_2$  normalization on the depth dimension of each tensor as in [2]. We reshape the normalized tensors into four  $d \times T$  matrices, where  $T = 14 \times 14$ .

Next, attention on the four layers is created from  $s_Q$ , the representation of the whole question defined above. We use a two-layer neural network having 724 hidden units with ReLU non-linearity to project  $s_Q$  to the scores of the four layers as

$$[s_1, s_2, s_3, s_4] = MLP(s_Q),$$
 (3)

which are then normalized by softmax to obtain four attention weights  $\alpha_1, \ldots, \alpha_4$ . The weighted sum of the above four matrices is computed, yielding a  $d \times T$  matrix  $V = [v_1, \ldots, v_T]$ , which is our representation of the input image. It stores the image feature at the t-th image region in its t-th column vector of size d.

### 3.2. Dense Co-Attention Layer

#### 3.2.1 Overview of the Architecture

We now describe the proposed dense co-attention layer; see Fig.2. It takes the question and image representations Q and V as inputs and then outputs their updated versions. We denote the inputs to the (l+1)-st layer by  $Q_l = [q_{l1},...,q_{lN}] \in \mathbb{R}^{d\times N}$  and  $V_l = [v_{l1},...,v_{lT}] \in \mathbb{R}^{d\times T}$ . For the first layer inputs, we set  $Q_0 = Q = [q_1,...,q_N]$  and  $V_0 = V = [v_1,...,v_T]$ .

The proposed architecture has the following properties. First, it is a co-attention mechanism [21]. Second, the coattention is *dense* in the sense that it considers every interaction between any word and any region. To be specific, our mechanism creates one attention map on regions per each word and creates one attention map on words per each region (see Fig.3). Third, it can be stacked as shown in Fig.1.

#### 3.2.2 Dense Co-attention Mechanism

**Basic method for attention creation** For the sake of explanation, we first explain the basic method for creation of attention maps, which we will extend later. Given  $Q_l$  and

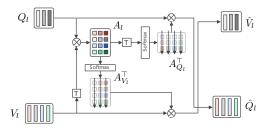


Figure 3: Computation of dense co-attention maps and attended representations of the image and question.

 $V_l$ , two attention maps are created as shown in Fig.3. Their computation starts with the affinity matrix

$$A_l = V_l^\top W_l Q_l, \tag{4}$$

where  $W_l$  is a learnable weight matrix. We normalize  $A_l$  in row-wise to derive attention maps on question words conditioned by each image region as

$$A_{Q_l} = \operatorname{softmax}(A_l), \tag{5}$$

and also normalize  $A_l$  in column-wise to derive attention maps on image regions conditioned by each question word as

$$A_{V_l} = \operatorname{softmax}(A_l^{\top}). \tag{6}$$

Note that each row of  $A_{Q_l}$  and  $A_{V_l}$  contains a single attention map.

Nowhere-to-attend and memory It often occurs at the creation and application of each attention map that there is no particular region or word that the model should attend. To deal with such cases, we add K elements to N question words as well as to T image regions, as in [30]. In [30], the authors only use K=1, but we found it effective to use K>1, which is expected to additionally serve as a memory for storing useful information [9]. More specifically, incorporating two matrices  $M_{Q_l} \equiv [q_{l \oslash 1}, ..., q_{l \oslash K}] \in \mathbb{R}^{d \times K}$  and  $M_{V_l} \equiv [v_{l \oslash 1}, ..., v_{l \oslash K}] \in \mathbb{R}^{d \times K}$ , which are learnable parameters, we augment the matrix  $Q_l$  and  $V_l$  in the row direction as  $\tilde{Q}_l = [q_{l1}, ..., q_{lN}, q_{l \oslash 1}, ..., q_{l \oslash K}] \in \mathbb{R}^{d \times (N+K)}$  and  $\tilde{V}_l = [v_{l1}, ..., v_{lT}, v_{l \oslash 1}, ..., v_{l \oslash K}] \in \mathbb{R}^{d \times (N+K)}$ . This augmentation of  $Q_l$  and  $V_l$  provides  $A_l$  of size  $(T+K) \times (N+K)$ ;  $A_{Q_l}$  and  $A_{V_l}$  are of size  $(T+K) \times (N+K)$  and  $(N+K) \times (T+K)$ , respectively.

**Parallel attention** In several studies [14, 29], multiple attention maps are created and applied to target features in a parallel manner, which provides multiple attended features, and then they are *fused by concatenation*. In [29], features are first linearly projected to multiple lower-dimensional spaces, for each of which the above attention function is performed. We adopt a similar approach that uses multiple

attention maps here, but we use *average* instead of concatenation for fusion of the multiple attended features, because we found it works better in our case.

To be specific, we linearly project the d-dimensional features (stored in the columns) of  $\tilde{V}_l$  and  $\tilde{Q}_l$  to multiple lower dimensional spaces. Let h be the number of lower dimensional spaces and  $d_h (\equiv d/h)$  be their dimension. We denote the linear projections by  $W_{\tilde{V}_l}^{(i)} \in \mathbb{R}^{d_h \times d}$  and  $W_{\tilde{Q}_l}^{(i)} \in \mathbb{R}^{d_h \times d}$   $(i=1,\ldots,h)$ . Then the affinity matrix between the projected features in the i-th space is given as

$$A_l^{(i)} = (W_{\tilde{V}_l}^{(i)} \tilde{V}_l)^\top (W_{\tilde{Q}_l}^{(i)} \tilde{Q}_l). \tag{7}$$

Attention maps are created from each affinity matrix by column-wise and row-wise normalization as

$$A_{Q_l}^{(i)} = \operatorname{softmax}\left(\frac{A_l^{(i)}}{\sqrt{d_h}}\right),\tag{8}$$

$$A_{V_l}^{(i)} = \operatorname{softmax}\left(\frac{A_l^{(i)\top}}{\sqrt{d_h}}\right). \tag{9}$$

As we employ multiplicative (or dot-product) attention as explained below, average fusion of multiple attended features is equivalent to averaging our attention maps as

$$A_{Q_l} = \frac{1}{h} \sum_{i=1}^{h} A_{Q_l}^{(i)}, \tag{10}$$

$$A_{V_l} = \frac{1}{h} \sum_{i=1}^{h} A_{V_l}^{(i)}.$$
 (11)

**Attended feature representations** We employ multiplicative attention to derive attended feature representations  $\hat{Q}_l$  and  $\hat{V}_l$  of the question and image, as shown in Fig.3. As  $A_{Q_l}$  and  $A_{V_l}$  store attention maps in their rows and their last K rows correspond to "nowhere-to-attend" or memory, we discard them when applying them to  $\tilde{Q}_l$  and  $\tilde{V}_l$  as

$$\hat{Q}_l = \tilde{Q}_l A_{Q_l} [\mathbf{1} : \mathbf{T}, :]^\top, \tag{12}$$

and

$$\hat{V}_l = \tilde{V}_l A_{V_l} [\mathbf{1} : \mathbb{N}, :]^\top, \tag{13}$$

respectively<sup>1</sup>. Note that  $\hat{Q}_l$  is the same size as  $V_l$  (i.e.  $d \times T$ ) and  $\hat{V}_l$  is the same size as  $Q_l$  (i.e.  $d \times N$ ).

# 3.2.3 Fusing Image and Question Representations

After computing the attended feature representations  $\hat{Q}_l$  and  $\hat{V}_l$ , we fuse the image and question representations, as shown in the right half of Fig.2. The matrix  $\hat{V}_l$  stores in

 $<sup>^1{\</sup>rm The}$  notation (1:T,:) indicates the submatrix in the first T rows, as in Python.

its n-th column the attended representation of the entire image conditioned on the n-th question word. Then, the n-th column vector  $\hat{v}_{ln}$  is fused with the representation  $q_{ln}$  of n-th question word by concatenation to form 2d-vector  $[q_{ln}^\top, \hat{v}_{ln}^\top]^\top$ . This concatenated vector is projected back to a d-dimensional space by a single layer network followed by the ReLU activation and residual connection as

$$q_{(l+1)n} = \text{ReLU}\left(W_{Q_l} \begin{bmatrix} q_{ln} \\ \hat{v}_{ln} \end{bmatrix} + b_{Q_l} \right) + q_{ln}, \quad (14)$$

where  $W_{Q_l} \in \mathbb{R}^{d \times 2d}$  and  $b_{Q_l} \in \mathbb{R}^d$  are learnable weights and biases. An identical network (with the same weights and biases) is applied to each question word  $(n=1,\ldots,N)$  independently, yielding  $Q_{l+1} = [q_{(l+1)1},\ldots,q_{(l+1)N}] \in \mathbb{R}^{d \times N}$ .

Similarly, the representation  $v_{lt}$  of t-th image region is concatenated with the representation  $\hat{q}_{lt}$  of the whole question words conditioned on the t-th image region, and then projected back to a d-dimensional space as

$$v_{(l+1)t} = \text{ReLU}\left(W_{V_l} \begin{bmatrix} v_{lt} \\ \hat{q}_{lt} \end{bmatrix} + b_{V_l} \right) + v_{lt}, \quad (15)$$

where  $W_{V_l} \in \mathbb{R}^{d \times 2d}$  and  $b_{V_l} \in \mathbb{R}^d$  are weights and biases. The application of an identical network to each image region  $(t=1,\ldots,T)$  yields  $V_{l+1}=[v_{(l+1)1},\ldots,v_{(l+1)T}] \in \mathbb{R}^{d \times T}$ .

It should be noted that the above two fully-connected networks have different parameters (i.e.,  $W_{Q_l}$ ,  $W_{V_l}$  etc.) for each layer l.

# 3.3. Answer Prediction

Given the final outputs  $Q_L$  and  $V_L$  of the last dense coattention layer, we predict answers. As they contain the representation of N question words and T image regions, we first perform self-attention function on each of them to obtain aggregated representations of the whole question and image. This is done for  $Q_L$  as follows: i) compute 'scores'  $s_{q_{L1}},\ldots,s_{q_{LN}}$  of  $q_{L1},\ldots,q_{LN}$  by applying an identical two-layer MLP with ReLU nonlinearity in its hidden layer; ii) then apply softmax to them to derive attention weights  $\alpha_1^Q,\ldots,\alpha_N^Q$ ; and iii) compute an aggregated representation by  $s_{Q_L}=\sum_{n=1}^N\alpha_n^Qq_{Ln}$ . Following the same procedure with an MLP with different weights, we derive attention weights  $\alpha_1^V,\ldots,\alpha_T^V$  and then compute an aggregated representation  $s_{V_L}$  from  $v_{L1},\ldots,v_{LT}$ .

Using  $s_{Q_L}$  and  $s_{V_L}$  thus computed, we predict answers. We consider three methods to do this here. The first one is to compute inner product between the sum of  $s_{Q_L}$  and  $s_{V_L}$  and  $s_{A}$ , the answer representation defined in Sec.3.1.1, as

(score of the answer encoded as  $s_A$ )

$$= \sigma \left( s_A^\top W \left( s_{Q_L} + s_{V_L} \right) \right), \quad (16)$$

where  $\sigma$  is the logistic function and W is a learnable weight matrix. The second and third ones are to use a MLP for computing scores for a set of predefined answers, which is a widely used approach in recent studies. The two differ in how to fuse  $s_{Q_L}$  and  $s_{V_L}$ , i.e., summation

$$(\text{score of answers}) = \sigma \Big( \text{MLP} \big( s_{Q_L} + s_{V_L} \big) \Big), \qquad (17)$$

or concatenation

(score of answers) = 
$$\sigma\left(\text{MLP}\left(\begin{bmatrix} s_{Q_L} \\ s_{V_L} \end{bmatrix}\right)\right)$$
, (18)

where MLP is a two layer MLP having 1024 hidden units with ReLU non-linearity. The first one is the most flexible, as it allows us to deal with any answers that are not considered at the time of training the entire network.

# 4. Experiments

In this section, we present results of the experiments conducted to evaluate the proposed method.

### 4.1. Datasets

We used two most popular datasets, VQA [2] and VQA 2.0 [8], for our experiments. VQA (also known as VQA 1.0) contains human-annotated question-answer pairs on 204,721 images from Microsoft COCO dataset [20]. There are three predefined splits of questions, *train*, *val* and *test* or *test-standard*, which consist of 248,349, 121,512, and 244,302 questions, respectively. There is also a 25% subset of the test-standard set referred to as *test-dev*. All the questions are categorized into three types: yes/no, number, and other. Each question has 10 free-response answers. VQA 2.0 is an updated version of VQA 1.0 and is the largest as of now. Compared with VQA 1.0, it contains more samples (443,757 *train*, 214,354 *val*, and 447,793 *test* questions), and is more balanced in term of language bias. We evaluate our models on the challenging Open-Ended task of both datasets.

As in [28], we choose correct answers appearing more than 5 times for VQA and 8 times for VQA 2.0 to form the set of candidate answers. Following previous studies, we train our network on train + val splits and report the test-dev and test-standard results from the VQA evaluation server (except for the ablation test shown below). We use the evaluation protocol of [2] in all the experiments.

# 4.2. Experimental Setup

For both of the datasets, we use the Adam optimizer with the parameters  $\alpha=0.001$ ,  $\beta_1=0.9$ , and  $\beta_2=0.99$ . During the training procedure, we make the learning rate  $(\alpha)$  decay at every 4 epochs for VQA and 7 epochs for VQA 2.0 with an exponential rate of 0.5. All models are trained up to



Figure 4: Typical examples of attended image regions and question words for complementary image-question pairs from VQA 2.0 dataset. Each row contains visualization for two pairs of the same question but different images and answers. The original image and question are shown along with their attention maps generated in the answer prediction layer. The brightness of image pixels and redness of words indicate the attention weights.

Table 1: Ablation study on each module of DCNs using the validation set of the Open-Ended task (VQA 2.0). \* indicates modules employed in the final model.

Category	Detail	Accuracy
Attention direction	$I \leftarrow Q$	60.95
	$I \rightarrow Q$	62.63
	$I \leftrightarrow Q^*$	62.94
Memory size $(K)$	1	62.53
	3*	62.94
	5	62.83
Number (h) of	2	62.82
parallel attention	4*	62.94
maps	8	62.81
Number $(L)$ of	1	62.43
stacked layers	2	62.82
	3*	62.94
	4	62.67
Attention in answer	Attention used*	62.94
prediction layer	Avg of features	61.63
Attention in image	Attention used*	62.94
extraction layer	Only last conv layer	62.39

16 and 21 epochs on VQA and VQA 2.0, respectively. To prevent overfitting, dropouts [27] are used after each fully connected layers with a dropout ratio p=0.3 and after the LSTM with a dropout ratio p=0.1. The batch size is set to 160 and 320 for VQA and VQA 2.0. We set the dimension d of the feature space in the dense co-attention layers (equivalently, the size of its hidden layers) to be 1024.

### 4.3. Ablation Study

The architecture of the proposed DCN is composed of multiple modules. To evaluate the contribution of each module to final prediction accuracy, we conducted ablation tests. Using VQA 2.0, we evaluated several versions of DCNs with different parameters and settings by training them on the train split and calculating its performance on the val split. The results are shown in Table 1.

The first block of the table shows the effects of image-question co-attention. The numbers are performances obtained by a DCN with only question-guided attention on images  $(I \leftarrow Q)$ , with only image-guided attention on question words  $(I \rightarrow Q)$ , and the standard DCN with co-attention  $(I \leftrightarrow Q)$ . The single-direction variants generates only attention in either side of the two paths in the dense co-attention layer; the rest of the computations remain the same. The network with co-attention performs the best, ver-

Table 2: Results of the proposed method along with published results of others on VQA 1.0 in similar conditions (i.e., a single model; trained without an external dataset).

Model	Test-dev				Test-standard			
	Overall	Other	Number	Yes/No	Overall	Other	Number	Yes/No
VQA team [2]	57.75	43.08	36.77	80.50	58.16	43.73	36.53	80.569
SMem [31]	57.99	43.12	37.32	80.87	58.24	43.48	37.53	80.80
SAN [32]	58.70	46.10	36.60	79.30	58.90	-	-	-
FDA [12]	59.24	45.77	36.16	81.14	59.54	-	-	-
DNMN [1]	59.40	45.50	38.60	81.10	59.40	-	-	-
HieCoAtt [21]	61.00	51.70	38.70	79.70	62.10	-	-	-
RAU [24]	63.30	53.00	39.00	81.90	63.20	52.80	38.20	81.70
DAN [23]	64.30	53.90	39.10	83.00	64.20	54.00	38.10	82.80
Strong Baseline [14]	64.50	55.20	39.10	82.20	64.60	55.20	39.10	82.00
MCB [6]	64.70	55.60	37.60	82.50	-	-	-	-
N2NMNs [11]	64.90	-	-	-	-	-	-	-
MLAN [35]	64.60	53.70	40.20	83.80	64.80	53.70	40.90	83.70
MLB [16]	65.08	54.87	38.21	84.14	65.07	54.77	37.90	84.02
MFB [36]	65.90	56.20	39.80	84.00	65.80	56.30	38.90	83.80
MF-SIG-T3 [5]	66.00	56.37	39.34	84.33	65.88	55.89	38.94	84.42
DCN (16)	66.43	56.23	42.37	84.75	66.39	56.23	41.81	84.53
DCN (17)	66.89	57.31	42.35	84.61	67.02	56.98	42.34	85.04
DCN (18)	66.83	57.44	41.66	84.48	66.66	56.83	41.27	84.61

Table 3: Results of the proposed method along with published results of others on VQA 2.0 in similar conditions (i.e., a single model; trained without an external dataset). DCN(number) indicates the DCN equipped with the prediction layer that uses equation (number) for score computation. \*: trained with external datasets. ‡: the winner of VQA challenge 2017, unpublished.

Model	Test-dev				Test-standard			
	Overall	Other	Number	Yes/No	Overall	Other	Number	Yes/No
VQA team-Prior [8]	-	-	-	-	25.98	01.17	00.36	61.20
VQA team-Language only [8]	-	-	-	-	44.26	27.37	31.55	67.01
VQA team-LSTM+CNN [8]	-	-	-	-	54.22	41.83	35.18	73.46
MCB [6] reported in [8]	-	-	-	-	62.27	53.36	38.28	78.82
MF-SIG-T3 * [5]	64.73	55.55	42.99	81.29	-	-	-	-
Adelaide Model * ‡ [28]	62.07	52.62	39.46	79.20	62.27	52.59	39.77	79.32
Adelaide + Detector * ‡ [28]	65.32	56.05	44.21	81.82	65.67	56.26	43.90	82.20
DCN (16)	66.87	57.26	46.61	83.51	66.97	57.09	46.98	83.59
DCN (17)	66.72	56.77	46.65	83.70	67.04	56.95	47.19	83.85
DCN (18)	66.60	56.72	46.60	83.50	67.00	56.90	46.93	83.89

ifying the effectiveness of our co-attention implementation.

The second block of the table shows the impacts of K, which is the row size of  $M_{Q_l}$  and  $M_{V_l}$  that are used for augmenting  $Q_l$  and  $V_l$ . This augmentation is originally introduced to be able to deal with "nowhere to attend", which can be implemented by K=1 [30]. However, we found that the use of K>1 improves performance to a certain extent, which we think is because  $M_{Q_l}$  and  $M_{V_l}$  work as external memory that can be used through attention mech-

anism [9]. As shown in the table, K=3 yields the best performance.

The third and fourth blocks of the table show choices of the number h of parallel attention maps and L of stacked layers. The best result was obtained for h=4 and L=3.

The last two blocks of the table show effects of the use of attention in the answer prediction layer and the image extraction layer; the use of attention improves accuracy by about 1.3% and 0.5%, respectively.

# 4.4. Comparison with Existing Methods

Table 2 shows the performance of our method on VQA 1.0 along with published results of others. The entries 'DCN (n)' indicate which method for score computation is employed from (16)-(18). It is seen from the table that our method outperforms the best published result (MF-SIG-T3) by a large margin of  $0.9\% \sim 1.1\%$  on both test-dev and test-standard sets. Furthermore, the improvements can be seen in all of the entries (*Other* with 1.1%, *Number* with 3.4%, *Yes/No* with 0.6% on test-standard set) implying the capacity of DCNs to model multiple types of complex relations between question-image pairs. Notably, we achieve significant improvements of 3.0% and 3.4% for the question type *Number* on test-dev and test-standard sets, respectively.

Table 2 also shows the performances of DCNs with a different answer prediction layer that uses (16), (17), and (18) for score computation. It is seen that (17) shows at least comparable performance to the others and even attains the best performance of 67.02% in test-standard set.

Table 3 shows comparisons of our method to previous published results on VQA 2.0 and also that of the winner of VQA 2.0 Challenge 2017 in both test-dev and test-standard sets. It is observed in Table 3 that our approach outperforms the state-of-the-art published method (MF-SIG-T3) by a large margin of 2.1% on test-dev set, even though the MF-SIG-T3 model was trained with VQA 2.0 and an augmented dataset (Visual Genome [19]). It is noted that the improvements are seen in all the question types (Other with 1.71%, *Number* with 3.66%, and *Yes/No* with 2.41%). Comparing our DCN with the winner of VOA 2.0 Challenge 2017, Adelaide model. Our best DCN (17) delivers 1.5% and 1.37% improvements in every question types over the Adelaide+Detector on test-dev and test-standard, respectively. It is worth to point out that the winner method uses a detector (Region Proposal Network) trained on annotated regions of Visual Genome dataset [19] to extract visual features; and that the model is trained using also an external dataset, i.e., the Visual Genome question answering dataset.

It should also be noted that while achieving the best performance in VQA dataset, the size of the DCNs (i.e., the number of parameters) is comparable or even smaller than the former state-of-the-art methods, as shown in Table 4.

# 4.5. Qualitative Evaluation

Complementary image-question pairs are available in VQA 2.0 [8], which are pairs of the same question and different images with different answers. To understand the behaviour of the trained DCN, we visualize attention maps that the DCN generates for some of the complementary image-question pairs. Specifically, we show multiplication of an input image and question with their attention maps  $\alpha_1^V,\ldots,\alpha_T^V$  and  $\alpha_1^Q,\ldots,\alpha_N^Q$  (defined in Sec.3.3) generated in the answer prediction layer. A typical example is shown

Table 4: Model sizes of DCNs and several bilinear fusion methods. The numbers include the parameters of LSTM networks and exclude those of ResNets.

Model	No. params
MCB [6]	63M
MLB [16]	25M
MFB [36]	46M
DCN (18)	32M
DCN (17)	31M
DCN (16)	28M

in Fig.4. Each row shows the results for two pairs of the same question and different images, from which we can observe that the DCN is able to look at right regions to find the correct answers. Then, the first column shows the results for two pairs of the same image and different questions. It is observed that the DCN focuses on relevant image regions and question words to produce answers correctly. More visualization results including failure cases are provided in the supplementary material.

### 5. Conclusion

In this paper, we present a novel network architecture for VQA named the dense co-attention network. The core of the network is the dense co-attention layer, which is designed to enable improved fusion of visual and language representations by considering dense symmetric interactions between the input image and question. The layer can be stacked to perform multi-step image-question interactions. The layer stack combined with the initial feature extraction step and the final answer prediction layer, both of which have their own attention mechanisms, form the dense co-attention network. The experimental results on two datasets, VQA and VQA 2.0, confirm the effectiveness of the proposed architecture.

# Acknowledgement

This work was partly supported by JSPS KAKENHI Grant Number JP15H05919, JST CREST Grant Number JPMJCR14D1, Council for Science, Technology and Innovation (CSTI), Cross-ministerial Strategic Innovation Promotion Program (Infrastructure Maintenance, Renovation and Management), and the ImPACT Program Tough Robotics Challenge of the Council for Science, Technology, and Innovation (Cabinet Office, Government of Japan).

#### References

[1] J. Andreas, M. Rohrbach, T. Darrell, and D. Klein. Learning to compose neural networks for question answering. In *The* 

- Association for Computational Linguistics (HLT-NAACL), 2016. 7
- [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. 3, 5, 7
- [3] D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align and translate. In *Inter*national Conference on Learning Representations (ICLR), 2015.
- [4] H. Ben-younes, R. Cadène, M. Cord, and N. Thome. MU-TAN: multimodal tucker fusion for visual question answering. In *IEEE International Conference on Computer Vision* (ICCV), 2017. 12
- [5] Z. Chen, Z. Yanpeng, H. Shuaiyi, T. Kewei, and M. Yi. Structured attentions for visual question answering. In *International Conference on Computer Vision (ICCV)*, 2017. 1, 2, 7, 12
- [6] 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. In *Empirical Methods in Natural Language Processing (EMNLP)*, 2016. 1, 2, 7, 8
- [7] Y. Gao, O. Beijbom, N. Zhang, and T. Darrell. Compact bilinear pooling. In *International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2
- [8] Y. Goyal, T. Khot, D. Summers-Stay, D. Batra, and D. Parikh. Making the V in VQA matter: Elevating the role of image understanding in Visual Question Answering. In *International Conference on Computer Vision and Pattern* Recognition (CVPR), 2017. 5, 7, 8
- [9] A. Graves, G. Wayne, and I. Danihelka. Neural turing machines. arXiv preprint arXiv:1410.5401, 2014. 4, 7
- [10] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2, 3
- [11] R. Hu, J. Andreas, M. Rohrbach, T. Darrell, and K. Saenko. Learning to reason: End-to-end module networks for visual question answering. In *IEEE International Conference on Computer Vision (ICCV)*, 2017. 7
- [12] I. Ilievski, S. Yan, and J. Feng. A focused dynamic attention model for visual question answering. arXiv preprint arXiv:1604.01485, 2016. 2, 7
- [13] R. Jozefowicz, W. Zaremba, and I. Sutskever. An empirical exploration of recurrent network architectures. *Journal of Machine Learning Research*, 2015. 11
- [14] V. Kazemi and A. Elqursh. Show, ask, attend, and answer: A strong baseline for visual question answering. arXiv preprint arXiv:1704.03162, 2017. 2, 4, 7
- [15] J.-H. Kim, S.-W. Lee, D. Kwak, M.-O. Heo, J. Kim, J.-W. Ha, and B.-T. Zhang. Multimodal Residual Learning for Visual QA. In *International Conference on Neural Information Processing Systems (NIPS)*, 2016. 2, 12
- [16] J.-H. Kim, K.-W. On, W. Lim, J. Kim, J.-W. Ha, and B.-T. Zhang. Hadamard product for low-rank bilinear pooling. In *International Conference on Learning Representa*tions (ICLR), 2017. 1, 2, 7, 8, 12

- [17] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Rep*resentations (ICLR), 2015. 11
- [18] R. Kiros, Y. Zhu, R. Salakhutdinov, R. S. Zemel, A. Torralba, R. Urtasun, and S. Fidler. Skip-thought vectors. In *International Conference on Neural Information Processing* Systems (NIPS), 2015. 12
- [19] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma, M. S. Bernstein, and L. Fei-Fei. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International Journal of Computer Vision (IJCV)*, 2017. 8
- [20] T. Lin, M. Maire, S. J. Belongie, L. D. Bourdev, R. B. Girshick, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft COCO: common objects in context. In European Conference on Computer Vision (ECCV), 2014. 5
- [21] J. Lu, J. Yang, D. Batra, and D. Parikh. Hierarchical question-image co-attention for visual question answering. In *International Conference on Neural Information Process*ing Systems (NIPS), 2016. 1, 2, 3, 7
- [22] B. McCann, J. Bradbury, C. Xiong, and R. Socher. Learned in translation: Contextualized word vectors. arXiv preprint arXiv:1708.00107, 2017. 3, 12
- [23] H. Nam, J. Ha, and J. Kim. Dual attention networks for multimodal reasoning and matching. arXiv preprint arXiv:1611.00471, 2016. 2, 7
- [24] H. Noh and B. Han. Training recurrent answering units with joint loss minimization for VQA. arXiv preprint arXiv:1606.03647, 2016. 2, 7
- [25] J. Pennington, R. Socher, and C. D. Manning. Glove: Global vectors for word representation. In *Empirical Methods in Natural Language Processing (EMNLP)*, 2014. 3, 11
- [26] K. J. Shih, S. Singh, and D. Hoiem. Where to look: Focus regions for visual question answering. In *International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [27] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research (JMLR)*, 2014. 6
- [28] D. Teney, P. Anderson, X. He, and A. van den Hengel. Tips and tricks for visual question answering: Learnings from the 2017 challenge. arXiv preprint arXiv:1708.02711, 2017. 2, 5, 7, 11
- [29] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. arXiv preprint arXiv:1706.03762, 2017. 4
- [30] C. Xiong, V. Zhong, and R. Socher. Dynamic coattention networks for question answering. In *International Conference on Learning Representations (ICLR)*, 2017. 4, 7
- [31] H. Xu and K. Saenko. Ask, attend and answer: Exploring question-guided spatial attention for visual question answering. In *European Conference on Computer Vision (ECCV)*, 2016. 2, 7
- [32] Z. Yang, X. He, J. Gao, L. Deng, and A. J. Smola. Stacked attention networks for image question answering. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 1, 2, 7

- [33] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson. How transferable are features in deep neural networks? In *International Conference on Neural Information Processing Systems (NIPS)*, 2014. 2
- [34] J. Yosinski, J. Clune, A. Nguyen, T. Fuchs, and H. Lipson. Understanding neural networks through deep visualization. In *ICML Workshop on Deep Learning*, 2015. 3
- [35] D. Yu, J. Fu, Y. Rui, and T. Mei. Multi-level attention networks for visual question answering. In *International Conference on Computer Vision and Pattern Recognition* (CVPR), 2017. 1, 2, 7
- [36] Z. Yu, J. Yu, J. Fan, and D. Tao. Multi-modal factorized bilinear pooling with co-attention learning for visual question answering. In *International Conference on Computer Vision* (*ICCV*), 2017. 1, 2, 7, 8
- [37] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *European Conference on Computer Vision (ECCV)*, 2014. 3

This document contains: more details of our experimental setups (Sec.A), the evaluation of effects in the employment of Contextualized Word Vectors (Sec.B), more visualization of attention maps generated in the answer prediction layer including failure cases (Sec.C), and an analysis of the attention mechanism employed in the image feature extraction (Sec.D).

# A. More Details of the Experimental Setups

In our experiments, images and questions are preprocessed as follows. All the images were resized to  $448 \times 448$  before feeding into the CNN. All the questions were tokenized using Python Natural Language Toolkit (nltk). We used the vocabulary provided by the CommonCrawl-840B Glove model for English word vectors [25], and set out-of-vocabulary words to unk. As mentioned in the main paper, we chose the correct answer appearing more than 5 times (= 3,014 answers) for VQA 1.0, and 8 times (= 3,113 answers) for VQA 2.0 as in [28]. We capped the maximum length of questions at 14 words and then performed dynamic unrolling for each question to allow for questions of different lengths.

Throughout the experiments, we used three-layer DCNs, that is, DCNs with three dense co-attention layers (L=3). This number of layers were chosen based on our preliminary experiments. The Bi-LSTM was initialized following the recommendation in [13] and all the other parameters were initialized as suggested by Glorot *et al*. In the training procedure, the ADAM [17] optimizer was used to train our model for 16 and 21 epochs on VQA and VQA 2.0 with batch size of 160 and 320, respectively; weight decay with rate of 0.0001 was added. We used exponential decay to gradually decrease the learning rate as

$$\alpha_{step} = 0.5^{\frac{\text{epochs}}{\text{decay epochs}}} \alpha,$$

where the initial learning rate  $\alpha$  was set to  $\alpha = 0.001$ , and the decay epochs was set to 4 and 7 epochs for VQA and VQA 2.0 in turn; we set  $\beta_1 = 0.9$ , and  $\beta_2 = 0.99$ .

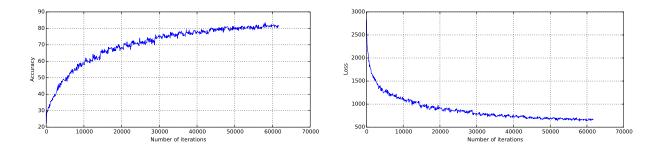


Figure 5: Learning curves for DCN.

# **B.** Effects of the Employment of Contextualized Word Vectors

To extract word features from input questions, some of the previous studies [15, 16, 5, 4] employed pretrained RNNs (specifically, GRU networks pre-trained with Skip-thought) [18]. In this study, we initially pursued a similar approach; we perform fine-tuning of a pretrained LSTM, specifically a two-layer Bi-LSTM trained as a CoVe (Context Vector) encoder [22]. Conducting comparative experiments, we eventually employ a single-layer Bi-LSTM with random initialization, as explained in the main paper. We report here the results of the experiments.

Table 5 shows the performances of DCNs with the CoVe-pretrained Bi-LSTM and with the randomly initialized Bi-LSTM. Note that the former is a two-layer model and the later has only one layer. Here, the VQA 2.0 test-dev dataset was used. It is observed that for DCNs with the answer prediction layer of (16), the one with the CoVe-pretrained model performs slightly better than the one with the randomly initialized model, but their differences are small. For DCNs with the answer prediction layers of (17) and (18), the one with the randomly initialized model performs better with a less number of parameters.

It should be noted, however, that the employment of CoVe-pretrained models, together with the answer prediction layer of (16), enables to compute meaningful answer representation  $(s_A)$  for answers that have not been seen before, i.e., those that are not included in training data. Table 6 shows the results of DCN (16) with the CoVe-pretrained model for *Multiple Choice* answers, which include a lot of unseen answers. This is not the case with DCNs (17) and (18) that compute scores of a fixed set of predetermined answers— the common approach of most of the recent studies.

Table 5: Performances of DCNs with the CoVe-pretrained LSTM and with the randomly initialized LSTM on the VQA 2.0 test-dev set.

Model	Overall	Other	Number	Yes/No	No. params
DCN (16) + CoVe	67.06	57.44	46.91	83.69	31M
DCN (16)	66.87	57.26	46.61	83.51	28M
DCN (17) + CoVe	66.21	56.71	46.01	82.72	34M
DCN (17)	66.72	56.77	46.65	83.70	31M
DCN (18) + CoVe	66.31	56.62	45.78	83.14	35M
DCN (18)	66.60	56.72	46.60	83.50	32M

Table 6: Effectiveness of DCN (16) + CoVe-pretrained LSTM on Multiple Choice answers.

Model	Test-dev				Test-std			
	Overall	Other	Number	Yes/No	Overall	Other	Number	Yes/No
DCN (16) + CoVe	71.37	66.10	45.48	84.39	71.20	65.93	44.13	84.23

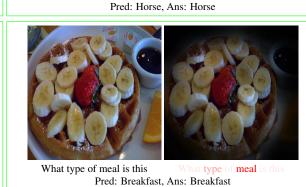
# C. Visualization of Attention Maps in the Answer Prediction Layer

We have shown a few examples of attention maps generated in the answer prediction layer of our DCNs in Fig.4 of the main paper. We show here more examples for success cases (Sec.C.1) and also for failure cases (Sec.C.2).

### C.1. Success Cases

We consider the visualization of complementary pairs to analyze the behaviour of our DCNs. Each row shows a complementary pair having the same question and different images. It can be seen from the examples shown below that the image and question attention maps are generated appropriately for most of success cases.





What is he sitting o

What is he sitting on



What type of meal is this

What type of meal is this





How many vases are in the photo Pred: 2, Ans: 2



How many vases are in the



How many vases are in the photo



How many vases are in the

Pred: 1, Ans: 1



of



What is the darker wall made What is the darker wall made

Pred: Brick, Ans: Brick





What is the darker wall made What is the darker wall made

Pred: Drywall, Ans: Drywall



What sport is this woman playing pred: Tennis, Ans: Tennis



What sport is this woman playing



What sport is this woman playing pl Pred: Frisbee, Ans: Frisbee



What **sport** is this woman playing

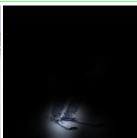


What color are the skiers shoes What color are the skiers shoes Pred: Yellow, Ans: Yellow





What color are the skiers shoes What color are the skiers shoes Pred: White, Ans: White





Does the man have a beard Pred: No, Ans: No



Does the man have a beard Pred: No, Ans: No



Is the sky blue or cloudy
Pred: Cloudy, Ans: Cloudy



Is the sky blue or cloudy
Pred: Blue, Ans: Blue



How many elephants How Pred: 2, Ans: 2

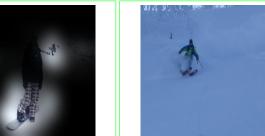


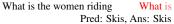
How many elephants

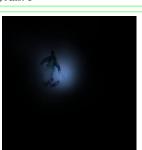




What is the women riding
Pred: Snowboard, Ans: Snowboard







What is the women riding

### C.2. Failure Cases

According to our analysis, failure cases can be categorized into the following four types:

- Type-1 Although the DCN is able to locate appropriate image regions and words, it fails to distinguish two different objects or concepts that have similar appearance. This may be attributable to that the extracted image features are not rich enough to distinguish them (e.g. *mutt* and *lab*; and *spoon* and *fork*).
- Type-2 Although the DCN is able to locate appropriate image regions and words, it fails to yield correct answers due to the bias of the dataset or missing instances of some objects/concepts in the dataset. For example, there are many samples of an *american flag* but no sample of a *dragon flag* in the training set.
- *Type-3* The DCN fails to locate appropriate image regions. This tends to occur when some image regions have similar appearance to the region that the DCN should attend, or the region that it should attend is too small.
- *Type-4* Although the DCN does yield conceptually correct answers, they are not listed in the given set of answers in the dataset and thus judged incorrect. For instance, while the given correct answer is *water*, the DCN outputs *beach*, which should also be correct, as in one of the examples below.

As in the above success cases, each row shows a complementary pair having the same question and different images. In each row, at least either one of the two has an erroneous prediction. The red bounding boxes indicate erroneous answers and the green ones indicate correct answers. The numbers in the failure examples indicate the error types we categorize above.



What breed of dog is this
Pred: Mutt, Ans: Lab (Error type: 1)



What breed of dog is this
Pred: Terrier, Ans: Terrier



What room is this
Pred: Bedroom, Ans: Bedroom



What room is this
Pred: Living room, Ans: Office (Error type: 1)



What is the name of the utensil
Pred: Fork, Ans: Fork



What is the name of the utensil What is the name of the utensil Pred: Fork, Ans: Spoon (*Error type: 1*)



How tall is he How tall is h
Pred: 5 feet, Ans: Tall (Error type: 1)

How tall is he
Pred: 5 feet, Ans: 6 feet (Error type: 2)









What is the color of pants the What is the color of pants the woman is wearing woman is wearing Pred: Plaid, Ans: Red and White (Error type: 4)





What is the color of pants the What is the color of pants the woman is wearing Pred: Green, Ans: Black (Error type: 4)





What color is lit up on the What color is lit up on the street lights Pred: Yellow, Ans: Green (Error type: 3)

Pred: Table, Ans: Plate (Error type: 4)





What color is lit up on the What color is lit up on the street lights Pred: White, Ans: None (Error type: 1)



Where is the fruit









Where is the fruit Where is the fruit Pred: Bowl, Ans: Bowl







How many tags are on the How many tags are on the suitcase Pred: 4, Ans: 3 (Error type: 1)





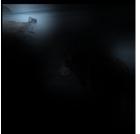
How many tags are on the How many tags are on the suitcase Pred: 0, Ans: 0



What landforms are behind the What landforms are behind the cows

Pred: Mountains, Ans: Mountains





What landforms are behind the What landforms are behind the cows

Pred: Beach, Ans: Water (Error type: 4)





What flag is that What flag is that Pred: American, Ans: American





What flag is that What flag is that Pred: American, Ans: Dragon (Error type: 2)





What is in the mug
Pred: Coffee, Ans: Coffee





What is in the mug
Pred: Wine, Ans: Butter (Error type: 1)





Where is this woman at Pred: Outside, Ans: Farmers market (*Error type: 4*)





Where is this woman at Pred: Market, Ans: Market

# D. Layer Attention in the Image Feature Extraction Step

As explained in the main paper (Sec.3.1), our DCN extracts visual features from an input image using a pre-trained ResNet at the initial step. The features are obtained by computing the weighted sum of the activations (i.e., outputs) of the four convolutional layers of the ResNet, where the attention weights generated conditioned on the input question are used. We examine here how this attention mechanism works for different types of questions. Specifically, utilizing the fifty five question types provided in the VQA-2.0, we compute the mean and standard deviation of the four attention weights for the questions belonging to each question type. We used all the questions in the validation set and our DCN trained only on train set for this computation.

Figure 6 shows the results. The bars in four colors represent the means of the four layer weights for each question type, and the thin black bars attached to the color bars represent their standard deviations. The fifty five question types are ordered by their similarity in the horizontal axis. From the plot, we can make the following observations:

- Layer 1 (the lowest one) has a certain level of weights only for *Yes/No* questions (shown on about the left half of the plots) and no weight for other types of questions (on the right half);
- Layer 2 has a small weight only for Yes/No questions and no weight for other types of questions;
- Layer 3 tends to have large weights for questions about colors (e.g., "what color") and questions about presence of a given object(s) (e.g., "are there" and "how many");
- Layer 4 (the highest one) has the largest attention weights in most of the question types, indicating its importance in answering them.
- Specific questions, such as "what color" and "what sport is", tend to have smaller standard deviations than nonspecific questions, such as "is the woman" and "do you".

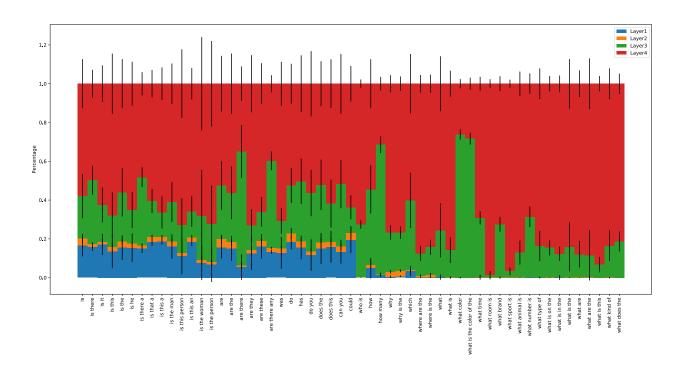


Figure 6: Statistics (means and standard deviations) of the attention weights on the four convolutional layers generated in the image feature extraction step for different types of questions.

# References

- [1] J. Andreas, M. Rohrbach, T. Darrell, and D. Klein. Learning to compose neural networks for question answering. In *The Association for Computational Linguistics (HLT-NAACL)*, 2016. 7
- [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. 3, 5, 7
- [3] D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align and translate. In *International Conference on Learning Representations (ICLR)*, 2015. 1
- [4] H. Ben-younes, R. Cadène, M. Cord, and N. Thome. MUTAN: multimodal tucker fusion for visual question answering. In *IEEE International Conference on Computer Vision (ICCV)*, 2017. 12
- [5] Z. Chen, Z. Yanpeng, H. Shuaiyi, T. Kewei, and M. Yi. Structured attentions for visual question answering. In *International Conference on Computer Vision (ICCV)*, 2017. 1, 2, 7, 12
- [6] 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. In *Empirical Methods in Natural Language Processing (EMNLP)*, 2016. 1, 2, 7, 8
- [7] Y. Gao, O. Beijbom, N. Zhang, and T. Darrell. Compact bilinear pooling. In *International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2
- [8] Y. Goyal, T. Khot, D. Summers-Stay, D. Batra, and D. Parikh. Making the V in VQA matter: Elevating the role of image understanding in Visual Question Answering. In *International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 5, 7, 8
- [9] A. Graves, G. Wayne, and I. Danihelka. Neural turing machines. arXiv preprint arXiv:1410.5401, 2014. 4, 7
- [10] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2, 3
- [11] R. Hu, J. Andreas, M. Rohrbach, T. Darrell, and K. Saenko. Learning to reason: End-to-end module networks for visual question answering. In *IEEE International Conference on Computer Vision (ICCV)*, 2017. 7
- [12] I. Ilievski, S. Yan, and J. Feng. A focused dynamic attention model for visual question answering. *arXiv preprint arXiv:1604.01485*, 2016. 2, 7
- [13] R. Jozefowicz, W. Zaremba, and I. Sutskever. An empirical exploration of recurrent network architectures. *Journal of Machine Learning Research*, 2015. 11

- [14] V. Kazemi and A. Elqursh. Show, ask, attend, and answer: A strong baseline for visual question answering. *arXiv* preprint *arXiv*:1704.03162, 2017. 2, 4, 7
- [15] J.-H. Kim, S.-W. Lee, D. Kwak, M.-O. Heo, J. Kim, J.-W. Ha, and B.-T. Zhang. Multimodal Residual Learning for Visual QA. In International Conference on Neural Information Processing Systems (NIPS), 2016. 2, 12
- [16] J.-H. Kim, K.-W. On, W. Lim, J. Kim, J.-W. Ha, and B.-T. Zhang. Hadamard product for low-rank bilinear pooling. In *International Conference on Learning Representations (ICLR)*, 2017. 1, 2, 7, 8, 12
- [17] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations* (ICLR), 2015. 11
- [18] R. Kiros, Y. Zhu, R. Salakhutdinov, R. S. Zemel, A. Torralba, R. Urtasun, and S. Fidler. Skip-thought vectors. In *International Conference on Neural Information Processing Systems (NIPS)*, 2015. 12
- [19] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma, M. S. Bernstein, and L. Fei-Fei. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International Journal of Computer Vision (IJCV)*, 2017. 8
- [20] T. Lin, M. Maire, S. J. Belongie, L. D. Bourdev, R. B. Girshick, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft COCO: common objects in context. In European Conference on Computer Vision (ECCV), 2014. 5
- [21] J. Lu, J. Yang, D. Batra, and D. Parikh. Hierarchical question-image co-attention for visual question answering. In *International Conference on Neural Information Processing Systems (NIPS)*, 2016. 1, 2, 3, 7
- [22] B. McCann, J. Bradbury, C. Xiong, and R. Socher. Learned in translation: Contextualized word vectors. *arXiv preprint* arXiv:1708.00107, 2017. 3, 12
- [23] H. Nam, J. Ha, and J. Kim. Dual attention networks for multimodal reasoning and matching. *arXiv preprint arXiv:1611.00471*, 2016. 2, 7
- [24] H. Noh and B. Han. Training recurrent answering units with joint loss minimization for VQA. arXiv preprint arXiv:1606.03647, 2016. 2, 7
- [25] J. Pennington, R. Socher, and C. D. Manning. Glove: Global vectors for word representation. In *Empirical Methods in Natural Language Processing (EMNLP)*, 2014. 3, 11
- [26] K. J. Shih, S. Singh, and D. Hoiem. Where to look: Focus regions for visual question answering. In *International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2
- [27] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research (JMLR)*, 2014. 6
- [28] D. Teney, P. Anderson, X. He, and A. van den Hengel. Tips and tricks for visual question answering: Learnings from the 2017 challenge. *arXiv preprint arXiv:1708.02711*, 2017. 2, 5, 7, 11
- [29] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. arXiv preprint arXiv:1706.03762, 2017. 4
- [30] C. Xiong, V. Zhong, and R. Socher. Dynamic coattention networks for question answering. In *International Conference on Learning Representations (ICLR)*, 2017. 4, 7
- [31] H. Xu and K. Saenko. Ask, attend and answer: Exploring question-guided spatial attention for visual question answering. In *European Conference on Computer Vision (ECCV)*, 2016. 2, 7
- [32] Z. Yang, X. He, J. Gao, L. Deng, and A. J. Smola. Stacked attention networks for image question answering. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 1, 2, 7
- [33] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson. How transferable are features in deep neural networks? In *International Conference on Neural Information Processing Systems (NIPS)*, 2014. 2
- [34] J. Yosinski, J. Clune, A. Nguyen, T. Fuchs, and H. Lipson. Understanding neural networks through deep visualization. In ICML Workshop on Deep Learning, 2015. 3
- [35] D. Yu, J. Fu, Y. Rui, and T. Mei. Multi-level attention networks for visual question answering. In *International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 1, 2, 7
- [36] Z. Yu, J. Yu, J. Fan, and D. Tao. Multi-modal factorized bilinear pooling with co-attention learning for visual question answering. In *International Conference on Computer Vision (ICCV)*, 2017. 1, 2, 7, 8
- [37] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *European Conference on Computer Vision (ECCV)*, 2014. 3