# Dual Recurrent Attention Units for Visual Question Answering

Ahmed Osman [1,2], Wojciech Samek [1]

[1] Fraunhofer Heinrich Hertz Institute, Berlin, Germany
[2] University of Freiburg, Freiburg, Germany

{ahmed.osman, wojciech.samek}@hhi.fraunhofer.de

## Abstract

*We propose an architecture for VQA which utilizes recurrent layers to generate visual and textual attention. The memory characteristic of the proposed recurrent attention units offers a rich joint embedding of visual and textual features and enables the model to reason relations between several parts of the image and question. Our single model outperforms the first place winner on the VQA 1.0 dataset, performs within margin to the current state-of-the-art ensemble model. We also experiment with replacing attention mechanisms in other state-of-the-art models with our implementation and show increased accuracy. In both cases, our recurrent attention mechanism improves performance in tasks requiring sequential or relational reasoning on the VQA dataset.*

## 1. Introduction

Although convolutional neural networks (CNNs) and recurrent neural networks (RNNs) have been successfully applied to various image and natural language processing tasks (cf. [1, 2, 3, 4]), these breakthroughs only slowly translate to multimodal tasks such as visual question answering (VQA) where the model needs to create a joint understanding of the image and question. Such multimodal tasks require joint visual and textual representations.

Since global features can hardly answer questions about certain local parts of the input, attention mechanisms have been extensively used in VQA recently [5, 6, 7, 8, 9, 10, 11, 12]. It attempts to make the model predict based on spatial or lingual context. However, most attention mechanisms used in VQA models are rather simple, consisting of two convolutional layers followed by a softmax to generate the attention weights which are summed over the image features. These shallow attention mechanisms may fail to select the relevant information from the joint representation of the question and image. Creating attention for complex questions, particularly sequential or relational reasoning questions, requires processing information in a sequen-
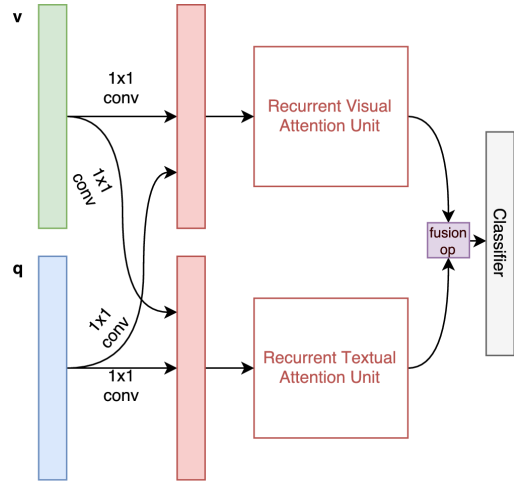


Figure 1. Diagram of the DRAU network.

tial manner which recurrent layers are better suited due to their ability to capture relevant information over an input sequence.

In this paper, we propose a RNN-based joint representation to generate visual and textual attention. We argue that embedding a RNN in the joint representation helps the model process information in a sequential manner and determine what is relevant to solve the task. We refer to the combination of RNN embedding and attention as Recurrent Textual Attention Unit (RTAU) and Recurrent Visual Attention Unit (RVAU) respective of their purpose. Furthermore, we employ these units in a fairly simple network, referred to as Dual Recurrent Attention Units (DRAU) network, and show improved results over several baselines. Finally, we enhance state-of-the-art models by replacing the model's default attention mechanism with RVAU.

Our main contributions are the following:

- We introduce a novel approach to generate soft attention. To the best of our knowledge, this is the first attempt to generate attention maps using recurrent neural networks. We provide quantitative and qualitative results showing performance improvements over the de-

fault attention used in most VQA models.

- Our attention modules are modular, thus, they can substitute existing attention mechanisms in most models fairly easily. We show that state-of-the-art models with RVAU "plugged-in" perform consistently better than their vanilla counterparts.

## 2. Related Work

This section discusses common methods that have been explored in the past for VQA.

**Bilinear representations**   Fukui et al. [7] use compact bilinear pooling to attend over the image features and combine it with the language representation. The basic concept behind compact bilinear pooling is approximating the outer product by randomly projecting the embeddings to a higher dimensional space using Count Sketch projection [13] and then exploiting Fast Fourier Transforms to compute an efficient convolution. An ensemble model using MCB won first place in VQA (1.0) 2016 challenge. Kim et al. [5] argues that compact bilinear pooling is still expensive to compute and shows that it can be replaced by element-wise product (Hadamard product) and a linear mapping (i.e. fully-connected layer) which gives a lower dimensional representation and also improves the model accuracy. Recently, Ben-younes et al. [14] proposed using Tucker decomposition [15] with a low-rank matrix constraint as a bilinear representation. They propose this fusion scheme in an architecture they refer to as MUTAN which as of this writing is the current state-of-the-art on the VQA 1.0 dataset.

**Attention-based**   Closely related to our work, Lu et al. [9] were the first to feature a co-attention mechanism that applies attention to both the question and image. Nam et al. [6] use a Dual Attention Network (DAN) that employs attention on both text and visual features iteratively to predict the result. The goal behind this is to allow the image and question attentions to iteratively guide each other in a synergistic manner.

**RNNs for VQA**   Using recurrent neural networks (RNNs) for VQA has been explored in the past. Xiong et al. [16] build upon the dynamic memory network from Kumar and Varaiya [17] and proposes DMN+. DMN+ uses episodic modules which contain attention-based Gated Recurrent Units (GRUs). Note that this is not the same as what we propose; Xiong et al. generate soft attention using convolutional layers and then uses it to substitute the update gate of the GRU. In contrast, our approach uses the recurrent layers to generate the attention. Noh and Han [8] propose recurrent answering units in which each unit is a complete module that can answer a question about an image. They use

joint loss minimization to train the units. However during testing, they use the first answering unit which was trained from other units through backpropagation.

**Notable mentions**   Kazemi and Elqursh [18] show that a simple model can get state-of-the-art results with proper training parameters. Wu et al. [19] construct a textual representation of the semantic content of an image and merges it with textual information sourced from a knowledge base. Ray et al. [20] introduce a task of identifying relevant questions for VQA. Kim et al. [21] apply residual learning techniques to VQA and propose a novel attention image attention visualization method using backpropagation.

## 3. Dual Recurrent Attention in VQA

We propose our method in this section. Figure 1 illustrates the flow of information in the DRAU model. Given an image and question, we create the input representations $v$ and $q$. Next, these features are combined by $1 \times 1$ convolutions into two separate branches. Then, the branches are passed to an RTAU and RVAU. Finally, the branches are combined using a fusion operation and fed to the final classifier. The full architecture of the network is depicted in Figure 2.

### 3.1. Input Representation

**Image representation**   We use the 152-layer "ResNet" pretrained CNN from He et al. [1] to extract image features. Similar to [7, 6], we resize the images to $448 \times 448$ and extract the last layer before the final pooling layer (res5c) with size $2048 \times 14 \times 14$. Finally, we use $l_2$ normalization on all dimensions. Recently, Anderson et al. [22] have shown that object-level features can provide a significant performance uplift compared to global-level features from pretrained CNNs. Therefore, we experiment with replacing the ResNet features with FRCNN [23] features with a fixed number of proposals per image ($K = 36$).

**Question representation**   We use a fairly similar representation as [7]. In short, the question is tokenized and encoded using an embedding layer followed by a tanh activation. We also exploit pretrained GloVe vectors [24] and concatenate them with the output of the embedding layer. The concatenated vector is fed to a two-layer unidirectional LSTM that contains 1024 hidden states each. In contrast to Fukui et al., we use all the hidden states of both LSTMs rather than concatenating the final states to represent the final question representation.

### 3.2. $1 \times 1$ **Convolution and PReLU**

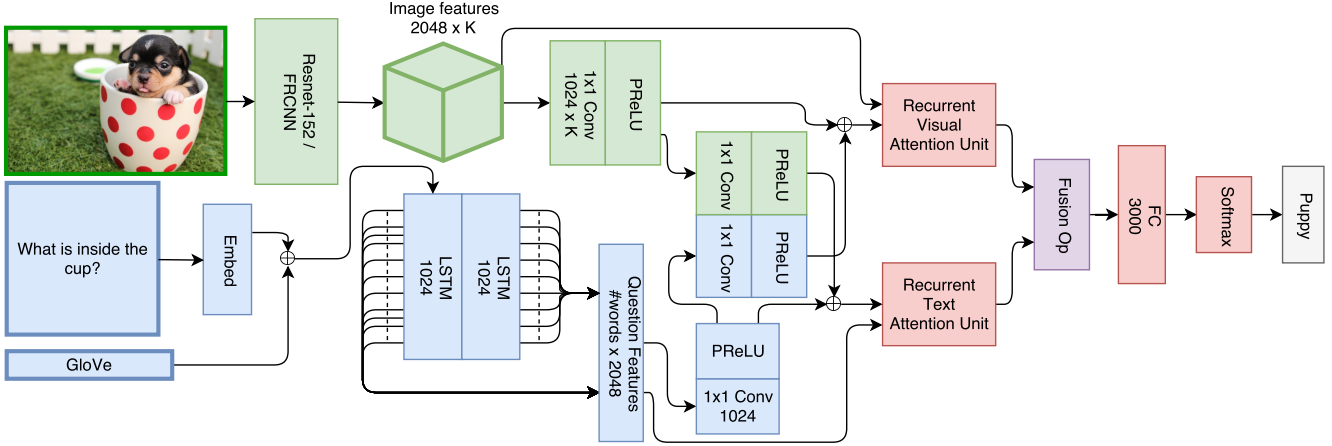We apply multiple $1 \times 1$ convolution layers in the network for mainly two reasons. First, they learn weights from

Figure 2. The proposed network. ⊕ denotes concatenation.

the image and question representations in the early layers. This is important especially for the image representation, since it was originally trained for a different task. Second, they are used to generate a common representation size. To obtain a joint representation, we apply $1 \times 1$ convolutions followed by PReLU activations [1] on both the image and question representations. Through empirical evidence, PReLU activations were found to reduce training time significantly and improve performance compared to ReLU and tanh activations. We provide these results in Section 4.

### 3.3. Recurrent Attention Units

The result from the above-mentioned layers is concatenated and fed to two separate recurrent attention units (RAU). Each RAU starts with another $1 \times 1$ convolution and PReLU activation:

$$c_a = \text{PReLU}\left(W_a\, x\right) \qquad (1)$$

where $W_a$ is the $1 \times 1$ convolution weights, $x$ is the input to the RAU, and $c_a$ is the output of the first PReLU.

Furthermore, we feed the previous output into an unidirectional LSTM:

$$h_{a,n} = \text{LSTM}\left(c_{a,n}\right) \qquad (2)$$

where $h_{a,n}$ is the hidden state at time $n$.

To generate the attention weights, we feed all the hidden states of the previous LSTM to a $1 \times 1$ convolution layer followed by a softmax function. The $1 \times 1$ convolution layer could be interpreted as the *number of glimpses* the model sees.

$$W_{att,n} = \text{softmax}\left(\text{PReLU}\left(W_g\, h_{a,n}\right)\right) \qquad (3)$$

where $W_g$ is the glimpses' weights and $W_{att,n}$ is the attention weight vector.

Next, we use the attention weights to compute a weighted average of the image and question features.

$$att_{a,n} = \sum_{n=1}^{N} W_{att,n}\, f_n \qquad (4)$$

where $f_n$ is the input representation and $att_{a,n}$ is the attention applied on the input. Finally, the attention maps are fed into a fully-connected layer followed by a PReLU activation. Figure 3 illustrates the structure of a RAU.

$$y_{att,n} = \text{PReLU}\left(W_{out}\, att_{a,n}\right) \qquad (5)$$

where $W_{out}$ is a weight vector of the fully connected layer and $y_{att,n}$ is the output of each RAU.

### 3.4. Reasoning layer

A fusion operation is used to merge the textual and visual branches. For DRAU, we experiment with using element-wise multiplication (Hadamard product) and MCB [7, 25]. The result of the fusion is given to a many-class classifier using the top 3000 frequent answers. We use a single-layer softmax with cross-entropy loss. This can be written as:

$$P_a = \text{softmax}\left(\text{fusion\_op}\left(y_{text}, y_{vis}\right) W_{ans}\right) \qquad (6)$$

where $y_{text}$ and $y_{vis}$ are the outputs of the RAUs, $W_{ans}$ represents the weights of the multi-way classifier, and $P_a$ is the probability of the top 3000 frequent answers.

The final answer $\hat{a}$ is chosen according to the following:

$$\hat{a} = \text{argmax}\, P_a \qquad (7)$$

## 4. Experiments and Results

Experiments are performed on the VQA 1.0 and 2.0 datasets [26, 27]. These datasets use images from the
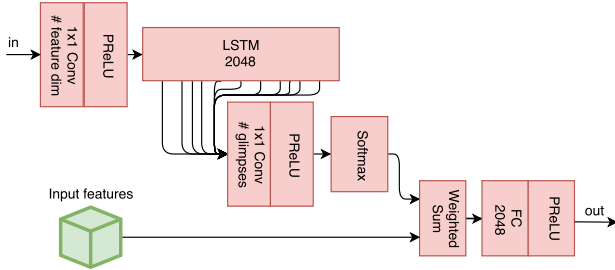
Figure 3. Recurrent Attention Unit.

| VQA 1.0 Validation Split | | | | |
|---|---|---|---|---|
| Open Ended Task | | | | |
| Baselines | Y/N | Num. | Other | All |
| Language only | 78.56 | 27.98 | 30.76 | 48.3 |
| Simple MCB | 78.64 | 32.98 | 39.79 | 54.82 |
| Joint LSTM | **79.90** | **36.96** | **49.58** | **59.34** |

Table 1. Evaluation of the baseline models on the VQA 1.0 *validation* split.

MS-COCO dataset [28] and generate questions and labels (10 labels per question) using Amazon's Mechanical Turk (AMT). Compared to VQA 1.0, VQA 2.0 adds more image-question pairs to balance the language prior present in the VQA 1.0 dataset. The ground truth answers in the VQA dataset are evaluated using human consensus.

$$\text{Acc}(a) = \min \left( \frac{\sum a \text{ is in human annotation}}{3}, 1 \right) \quad (8)$$

We evaluate our results on the *validation*, *test-dev*, *test-std* splits of each dataset. Models evaluated on the validation set use *train* and Visual Genome for training. For the other splits, we include the validation set in the training data. However, the models using FRCNN features do not use data augmentation with Visual Genome.

To train our model, we use Adam [29] for optimization with $\beta_1 = 0.9$, $\beta_2 = 0.999$, and an initial learning rate of $\epsilon = 7 \times 10^{-4}$. The final model is trained with a small batch size of 32 for 400K iterations. We did not fully explore tuning the batch size which explains the relatively high number of training iterations. Dropout ($p = 0.3$) is applied after each LSTM and after the fusion operation. All weights are initialized as described in [30] except LSTM layers which use an uniform weight distribution. The pre-trained ResNet was fixed during training due to the massive computational overhead of fine-tuning the network for the VQA task. While VQA datasets provide 10 answers per image-question pair, we sample one answer randomly for each training iteration.

## 4.1. VQA 1.0 Experiments

During early experiments, the VQA 2.0 dataset was not yet released. Thus, the baselines and early models were evaluated on the VQA 1.0 dataset. While building the final model, several parameters were changed, mainly, the learning rate, activation functions, dropout value, and other modifications which we discuss in this section.

**Baselines** We started by designing three baseline architectures. The first baseline produced predictions solely from the question while totally ignoring the image. The

model used the same question representation described in [7] and passed the output to a softmax 3000-way classification layer. The goal of this architecture was to assess the extent of the language bias present in VQA.

The second baseline is a simple joint representation of the image features and the language representation. The representations were combined using the compact bilinear pooling from [25]. We chose this method specifically because it was shown to be effective by Fukui et al. [7]. The main objective of this model is to measure how a robust pooling method of multimodal features would perform on its own without a deep architecture or attention. We refer to this model as *Simple MCB*.

For the last baseline, we substituted the compact bilinear pooling from Simple MCB with an LSTM consisting of hidden states equal to the image size. A $1 \times 1$ convolutional layer followed by a tanh activation were used on the image features prior to the LSTM, while the question representation was replicated to have a common embedding size for both representations This model is referred to as *Joint LSTM*.

We begin by testing our baseline models on the VQA 1.0 validation set. As shown in Table 1, the language-only baseline model managed to get 48.3% overall. More impressively, it scored 78.56% on Yes/No questions. The *Simple MCB* model further improves the overall performance, although little improvement is gained in the binary Yes/No tasks. Replacing MCB with our basic *Joint LSTM* embedding improves performance across the board.

**Modifications to the Joint LSTM Model** We test several variations of the *Joint LSTM* baseline which are highlighted in Table 2. Using PReLU activations has helped in two ways. First, it reduced time for convergence from 240K iterations to 120K. Second, the overall accuracy has improved, especially in the *Other* category. The next modifications were inspired by the results from [18]. We experimented with appending positional features which can be described as the coordinates of each pixel to the depth/feature dimension of the image representation. When unnormalized with respect to the other features, it worsened results significantly, dropping the overall accuracy by over 2 points.

| VQA 1.0 Validation Split Open Ended Task | | | | |
|---|---|---|---|---|
| Model | Y/N | Num. | Other | All |
| Joint LSTM baseline | **79.90** | **36.96** | 49.58 | 59.34 |
| PReLU | 79.61 | 36.21 | **50.77** | 59.74 |
| Pos. features | 79.68 | 36.52 | 46.59 | 57.71 |
| Pos. features (norm.) | 79.69 | 36.36 | 50.69 | **59.75** |
| High dropout | 79.03 | 34.84 | 47.25 | 57.59 |
| Extra FC | 78.86 | 33.51 | 45.57 | 56.51 |

Table 2. Evaluation of the Joint LSTM model and its modifications on the VQA 1.0 *validation* split.

Normalizing positional features did not have enough of a noticeable improvement (0.01 points overall) to warrant its effectiveness. Next, all dropout values are increased from 0.3 to 0.5 deteriorated the network's accuracy, particularly in the Number and Other categories. The final modification was inserting a fully connected layer with 1024 hidden units before the classifier, which surprisingly dropped the accuracy massively.

### 4.2. VQA 2.0 Experiments

After the release of VQA 2.0, we shifted our empirical evaluation towards the newer dataset. First, we retrain and retest our best performing VQA 1.0 model *Joint LSTM* as well as several improvements and modifications.

Since VQA 2.0 was built to reduce the language prior and bias inherent in VQA, the accuracy of *Joint LSTM* drops significantly as shown in Table 3. Note that all the models that were trained so far do not have explicit visual or textual attention implemented. Our first network with explicit visual attention, *RVAU*, shows an accuracy jump by almost 3 points compared to the Joint LSTM model. This result highlights the importance of attention for good performance in VQA. Training the *RVAU* network as a multi-label task ($RVAU_{multilabel}$), i.e. using all available annotations at each training iteration, drops the accuracy horribly. This is the biggest drop in performance so far. This might be caused by the variety of annotations in VQA for each question which makes the task for optimizing all answers at once much harder.

**DRAU Evaluation**   The addition of RTAU marks the creation of our *DRAU* network. The *DRAU* model shows favorable improvements over the *RVAU* model. Adding textual attention improves overall accuracy by 0.56 points. Substituting the PReLU activations with ReLU ($DRAU_{ReLU}$) massively drops performance. While further training might have helped the model improve, PReLU offers much faster

[1]Concurrent Work

| VQA 2.0 Validation Split Open Ended Task | | | | |
|---|---|---|---|---|
| Model | Y/N | Num. | Other | All |
| Joint LSTM w/PReLU | 72.04 | 37.95 | 48.58 | 56.00 |
| RVAU | 74.59 | 37.75 | 52.81 | 59.02 |
| $RVAU_{multilabel}$ | **77.53** | 36.05 | 40.18 | 53.67 |
| $DRAU_{Hadamard fusion}$ | 76.62 | **38.92** | 52.09 | 59.58 |
| $DRAU_{answer vocab = 5k}$ | 76.33 | 38.21 | 51.85 | 59.27 |
| $DRAU_{ReLU}$ | 72.69 | 34.92 | 45.05 | 54.11 |
| $DRAU_{no final dropout}$ | **77.02** | 38.26 | 50.17 | 58.69 |
| $DRAU_{high final dropout}$ | 76.47 | 38.71 | **52.52** | **59.71** |
| MCB [26] | - | - | - | 59.14 |
| Kazemi and Elqursh [18][1] | - | - | - | 59.67 |

Table 3. Evaluation of RVAU and DRAU-based models on the VQA 2.0 *validation* split.

convergence. Increasing the value of the dropout layer after the fusion operation ($DRAU_{high final dropout}$) improves performance by 0.13 points, in contrast to the results of the *Joint LSTM model* on VQA 1.0. Note that on the VQA 1.0 tests, we changed the values of all layers that we apply dropout on, but here we only change the last one after the fusion operation. Totally removing this dropout layer worsens accuracy. This suggests that the optimal dropout value should be tuned per-layer.

We test a few variations of *DRAU* on the test-dev set. We can observe that VQA benefits from more training data; the same *DRAU* network performs better (62.24% vs. 59.58%) thanks to the additional data. Most of the literature resize the original ResNet features from $224 \times 224$ to $448 \times 448$. To test the effect of this scaling, we train a *DRAU* variant with the original ResNet size ($DRAU_{small}$). Reducing the image feature size from $2048 \times 14 \times 14$ to $2048 \times 7 \times 7$ adversely affects accuracy as shown in Table 4. Adding more glimpses significantly reduces the model's accuracy ($DRAU_{glimpses = 4}$). A cause of this performance drop could be related to the fact that LSTMs process the input in a one-dimensional fashion and thus decide that each input is either relevant or non-relevant. This might explain why the attention maps of *DRAU* separate the objects from the background in two glimpses as we will mention in Section 5. 2D Grid LSTMs [31] might help remove this limitation. Removing the extra data from Visual Genome hurts the model's accuracy. That supports the fact that VQA is very diverse and that extra data helps the model perform better. Finally, substituting Hadamard product of MCB in the final fusion operation boosts the network's accuracy significantly by 1.17 points ($DRAU_{MCB fusion}$).

As mentioned in Section 3.1, we experiment replacing the global ResNet features with object-level features as sug-

| VQA 2.0 Test-Dev Split | | | | |
|---|---|---|---|---|
| Open Ended Task | | | | |
| Model | Y/N | Num. | Other | All |
| DRAU$_{\text{Hadamard fusion}}$ | 78.27 | 40.31 | 53.57 | 62.24 |
| DRAU$_{\text{small}}$ | 77.53 | 38.78 | 49.93 | 60.03 |
| DRAU$_{\text{glimpses} = 4}$ | 76.82 | 39.15 | 51.07 | 60.32 |
| DRAU$_{\text{no genome}}$ | 79.63 | 39.55 | 51.81 | 61.88 |
| DRAU$_{\text{MCB fusion}}$ | 78.97 | 40.06 | 55.47 | 63.41 |
| DRAU$_{\text{FRCNN features}}$ | **82.85** | **44.78** | **57.4** | **66.45** |

Table 4. Evaluation of later DRAU-based models on the VQA 2.0 *test-dev* split.

| VQA 2.0 Test-dev Split | | | | |
|---|---|---|---|---|
| Open Ended Task | | | | |
| Model | Y/N | Num. | Other | All |
| MCB [7][3] | 78.41 | 38.81 | 53.23 | 61.96 |
| MCB w/RVAU | 77.31 | **40.12** | **54.64** | 62.33 |
| MUTAN [14] | 79.06 | 38.95 | 53.46 | 62.36 |
| MUTAN w/RVAU | **79.33** | 39.48 | 53.28 | **62.45** |

Table 5. Results of state-of-the-art models with RVAU.

gested by [22]. This change provides a significant performance increase of 3.04 points (*DRAU$_{\text{FRCNN features}}$*).

### 4.3. Transplanting RVAU in other models

To verify the effectiveness of the recurrent attention units, we replace the attention layers in MCB and MUTAN [14] with RVAU.

For MCB [7] we remove all the layers after the first MCB operation until the first 2048-d output and replace them with RVAU. Due to GPU memory constraints, we reduced the size of each hidden unit in RVAU's LSTM from 2048 to 1024. In the same setting, RVAU significantly helps improve the original MCB model's accuracy as shown in Table 5. The most noticeable performance boost can be seen in the number category, which supports our hypothesis that recurrent layers are more suited for sequential reasoning.

Furthermore, we test RVAU in the MUTAN model [14]. The authors use a multimodal vector with dimension size of 510 for the joint representations. For coherence, we change the usual dimension size in RVAU to 510. At the time of this writing, the authors have not released results on VQA 2.0 using a single model rather than a model ensemble. Therefore, we train a single-model MUTAN using the authors' implementation.[2] The story does not change here, RVAU improves the model's overall accuracy.

### 4.4. DRAU versus the state-of-the-art

**VQA 1.0** Table 6 shows a comparison between DRAU and other state-of-the-art models. Excluding model ensembles, DRAU performs favorably against other models. To the best of our knowledge, [5] has the best single model performance of 65.07% on the *test-std* split which is very close our best model (65.03%). Small modifications or hyperparameter tuning could push our model further. Finally, the FRCNN image features boosts the model's performance close to the state-of-the-art ensemble model.

**VQA 2.0** Our model DRAU$_{\text{MCB fusion}}$ landed the 8th place in the VQA 2.0 Test-standard task.[4] Currently, all reported submissions that outperform our single model use model ensembles. Using FRCNN features boosted the model's performance to outperform some of the ensemble models (66.85%). The first place submission [22] reports using an ensemble of 30 models. In their report, the best single model that uses FRCNN features achieves 65.67% on the *test-standard* split which is outperformed by our best single model DRAU$_{\text{FRCNN features}}$.

## 5. DRAU versus MCB

In this section, we provide qualitative results that highlight the effect of the recurrent layers compared to the MCB model.

The strength of RAUs is notable in tasks that require sequentially processing the image or relational/multi-step reasoning. In the same setting, DRAU outperforms MCB in counting questions. This is validated in a subset of the validation split questions in the VQA 2.0 dataset as shown in Figure 4. Figure 5 shows some qualitative results between DRAU and MCB. For fair comparison we compare the first attention map of MCB with the second attention map of our model. We do so because the authors of MCB [7] visualize the first map in their work[5]. Furthermore, the first glimpse of our model seems to be the complement of the second attention, i.e. the model separates the background and the target object(s) into separate attention maps. We have not tested the visual effect of more than two glimpses on our model.

In Figure 5, it is clear that the recurrence helps the model attend to multiple targets as apparent in the difference of the attention maps between the two models. DRAU seems to also know how to count the right object(s). The top right example in Figure 5 illustrates that DRAU is not easily fooled by counting whatever object is present in the image but rather the object that is needed to answer the question. This

| VQA 1.0 Open Ended Task | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | Test-dev | | | | Test-standard | | | |
| Model | Y/N | Num. | Other | All | Y/N | Num. | Other | All |
| SAN [10] | 79.3 | 36.6 | 46.1 | 58.7 | - | - | - | 58.9 |
| DMN+ [16] | 80.5 | 36.8 | 48.3 | 60.3 | - | - | - | 60.4 |
| MRN [21] | 82.28 | 38.82 | 49.25 | 61.68 | 82.39 | 38.23 | 49.41 | 61.84 |
| HieCoAtt [9] | 79.7 | 38.7 | 51.7 | 61.8 | - | - | - | 62.1 |
| RAU [8] | 81.9 | 39.0 | 53.0 | 63.3 | 81.7 | 38.2 | 52.8 | 63.2 |
| DAN [6] | 83.0 | 39.1 | 53.9 | 64.3 | 82.8 | 38.1 | 54.0 | 64.2 |
| MCB [7] ($e = 7$) | 83.4 | 39.8 | 58.5 | 66.7 | 83.24 | 39.47 | 58.00 | 66.47 |
| MLB [5] (1 model) | - | - | - | - | 84.02 | 37.90 | 54.77 | 65.07 |
| MLB [5] ($e = 7$) | 84.57 | 39.21 | 57.81 | 66.77 | 84.61 | 39.07 | 57.79 | 66.89 |
| MUTAN [14] ($e = 5$) | 85.14 | 39.81 | 58.52 | 67.42 | 84.91 | 39.79 | 58.35 | 67.36 |
| $\text{DRAU}_{\text{Hadamard fusion}}$ | 82.73 | 38.18 | 54.43 | 64.3 | - | - | - | - |
| $\text{DRAU}_{\text{MCB fusion}}$ | 82.44 | 38.22 | 56.30 | 65.1 | 82.41 | 38.33 | 55.97 | 65.03 |
| $\text{DRAU}_{\text{FRCNN features}}$ | 84.92 | 39.16 | 57.70 | 66.86 | 84.87 | 40.02 | 57.91 | 67.16 |

Table 6. DRAU compared to the state-of-the-art on the VQA 1.0 dataset.$e = n$ corresponds to a model ensemble of size $n$.

| VQA 2.0 Open Ended Task | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | Test-dev | | | | Test-standard | | | |
| Model | Y/N | Num. | Other | All | Y/N | Num. | Other | All |
| UPC | 67.1 | 31.54 | 25.46 | 43.3 | 66.97 | 31.38 | 25.81 | 43.48 |
| MIC_TJ | 69.02 | 34.52 | 35.76 | 49.32 | 69.22 | 34.16 | 35.97 | 49.56 |
| neural-vqa-attention[10] | 70.1 | 35.39 | 47.32 | 55.35 | 69.77 | 35.65 | 47.18 | 55.28 |
| CRCV_REU | 73.91 | 36.82 | 54.85 | 60.65 | 74.08 | 36.43 | 54.84 | 60.81 |
| VQATeam_MCB [26] | 78.41 | 38.81 | 53.23 | 61.96 | 78.82 | 38.28 | 53.36 | 62.27 |
| DCD_ZJU[32] | 79.84 | 38.72 | 53.08 | 62.47 | 79.85 | 38.64 | 52.95 | 62.54 |
| VQAMachine [33] | 79.4 | 40.95 | 53.24 | 62.62 | 79.82 | 40.91 | 53.35 | 62.97 |
| POSTECH | 78.98 | 40.9 | 55.35 | 63.45 | 79.32 | 40.67 | 55.3 | 63.66 |
| UPMC-LIP6[14] | 81.96 | 41.62 | 57.07 | 65.57 | 82.07 | 41.06 | 57.12 | 65.71 |
| LV_NUS[34] | 81.95 | 48.31 | 59.99 | 67.71 | 81.92 | 48.38 | 59.63 | 67.64 |
| DLAIT | 82.94 | 47.08 | 59.94 | 67.95 | 83.17 | 46.66 | 60.15 | 68.22 |
| HDU-USYD-UNCC[35] | 84.39 | 45.76 | 59.14 | 68.02 | 84.5 | 45.39 | 59.01 | 68.09 |
| Adelaide-Teney ACRV MSR[36] | 85.24 | 48.19 | 59.88 | 69.00 | 85.54 | 47.45 | 59.82 | 69.13 |
| $\text{DRAU}_{\text{Hadamard fusion}}$ | 78.27 | 40.31 | 53.58 | 62.24 | 78.86 | 39.91 | 53.76 | 62.66 |
| $\text{DRAU}_{\text{MCB fusion}}$ | 78.97 | 40.06 | 55.48 | 63.41 | 79.27 | 40.15 | 55.55 | 63.71 |
| $\text{DRAU}_{\text{FRCNN features}}$ | 82.85 | 44.78 | 57.4 | 66.45 | 83.35 | 44.37 | 57.63 | 66.85 |

Table 7. DRAU compared to the current submissions on the VQA 2.0 dataset.

property also translates to questions that require relational reasoning. The second column in Figure 5 demonstrates how DRAU can attend the location required to answer the question based on the textual and visual attention maps.

## 6. Conclusion

We proposed an architecture for VQA with a novel attention unit, termed the Recurrent Attention Unit (RAU). The recurrent layers help guide the textual and visual attention since the network can reason relations between several parts of the image and question. We provided quantitative and qualitative results indicating the usefulness of a recurrent attention mechanism. Our DRAU model showed improved performance in tasks requiring sequential/complex reasoning such as counting or relational reasoning over [7], the winners of the VQA 2016 challenge. In VQA 1.0, we achieved near state-of-the-art results for single model performance with our DRAU network (65.03 vs. 65.07 [5]). Adding the FRCNN features gets the model within margin
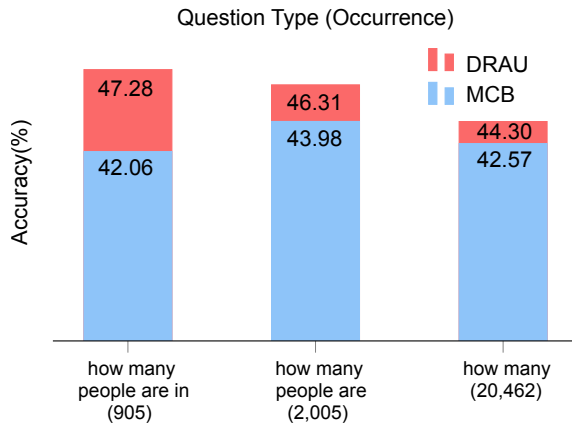
## Question Type (Occurrence)



Figure 4. Results on questions that require counting in the VQA 2.0 validation set.

of the state-of-the-art 5-model ensemble MUTAN [14]. Finally, we demonstrated that substituting the visual attention mechanism in other networks, MCB [7] and MUTAN [14], consistently improves their performance. In future work we will investigate implicit recurrent attention mechanism using recently proposed explanation methods [37, 38].

## References

[1] K. He, X. Zhang, S. Ren, and J. Sun, "Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification," in *IEEE International Conference on Computer Vision (ICCV)*, 2015, pp. 1026–1034.

[2] S. Bosse, D. Maniry, K.-R. Müller, T. Wiegand, and W. Samek, "Deep neural networks for no-reference and full-reference image quality assessment," *IEEE Transactions on Image Processing*, vol. 27, no. 1, pp. 206–219, 2018.

[3] D. Bahdanau, K. Cho, and Y. Bengio, "Neural Machine Translation by Jointly Learning to Align and Translate," in *International Conference on Representation Learning (ICLR)*, 2015.

[4] R. Nallapati, B. Zhou, C. Gulcehre, B. Xiang *et al.*, "Abstractive text summarization using sequence-to-sequence rnns and beyond," *arXiv:1602.06023*, 2016.

[5] J.-H. Kim, K.-W. On, J. Kim, J.-W. Ha, and B.-T. Zhang, "Hadamard Product for Low-rank Bilinear Pooling," in *International Conference on Representation Learning (ICLR)*, 2017.

[6] H. Nam, J.-W. Ha, and J. Kim, "Dual Attention Networks for Multimodal Reasoning and Matching," in *IEEE International Conference on Computer Vision (CVPR)*, 2017, pp. 299–307.

[7] A. Fukui, D. H. Park, D. Yang, A. Rohrbach, T. Darrell, and M. Rohrbach, "Multimodal Compact Bilinear Pooling for Visual Question Answering and Visual Grounding," in *Empirical Methods in Natural Language Processing (EMNLP)*, 2016, pp. 457–468.

[8] H. Noh and B. Han, "Training Recurrent Answering Units with Joint Loss Minimization for VQA," *arXiv:1606.03647*, 2016.

[9] J. Lu, J. Yang, D. Batra, and D. Parikh, "Hierarchical Question-Image Co-Attention for Visual Question Answering," in *Advances in Neural Information Processing Systems (NIPS)*, 2016, pp. 289–297.

[10] Z. Yang, X. He, J. Gao, L. Deng, and A. Smola, "Stacked Attention Networks for Image Question Answering," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 21–29.

[11] K. Chen, J. Wang, L.-C. Chen, H. Gao, W. Xu, and R. Nevatia, "ABC-CNN: An Attention Based Convolutional Neural Network for Visual Question Answering," *arXiv:1511.05960*, 2015.

[12] 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, pp. 451–466.

[13] M. Charikar, K. Chen, and M. Farach-Colton, "Finding frequent items in data streams," *Theoretical Computer Science*, vol. 312, no. 1, pp. 3–15, 2004.

[14] H. Ben-younes, R. Cadene, M. Cord, and N. Thome, "MUTAN: Multimodal Tucker Fusion for Visual Question Answering," *arXiv:1705.06676*, 2017.

[15] L. R. Tucker, "Some mathematical notes on three-mode factor analysis," *Psychometrika*, vol. 31, no. 3, pp. 279–311, 1966.

[16] C. Xiong, S. Merity, and R. Socher, "Dynamic Memory Networks for Visual and Textual Question Answering," in *International Conference on Machine Learning (ICML)*, 2016, pp. 2397–2406.

[17] P. R. Kumar and P. Varaiya, *Stochastic systems: Estimation, identification, and adaptive control*. SIAM, 2015.

[18] V. Kazemi and A. Elqursh, "Show, Ask, Attend, and Answer: A Strong Baseline For Visual Question Answering," *arXiv:1704.03162*, 2017.

[19] Q. Wu, P. Wang, C. Shen, A. Dick, and A. van den Hengel, "Ask Me Anything: Free-form Visual Question Answering Based on Knowledge from External Sources," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 4622–4630.

[20] A. Ray, G. Christie, M. Bansal, D. Batra, and D. Parikh, "Question Relevance in VQA: Identifying Non-Visual And False-Premise Questions," in *Empirical Methods in Natural Language Processing (EMNLP)*, 2016, pp. 919–924.
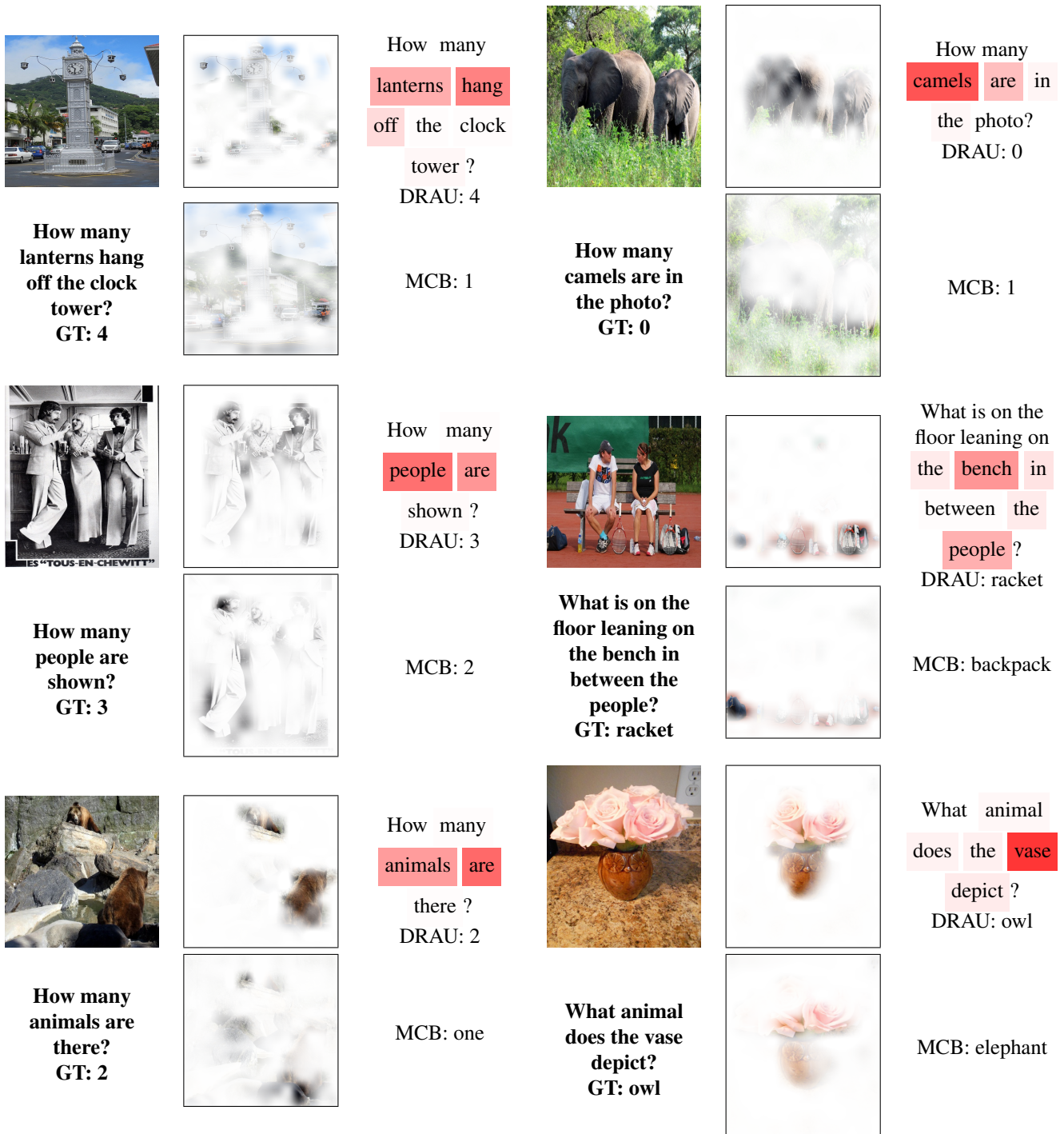
Figure 5. DRAU vs. MCB Qualitative examples. Attention maps for both models shown, only DRAU has textual attention.

[21] 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 *Advances in Neural Information Processing Systems 29*, 2016, pp. 361–369.

[22] P. Anderson, X. He, C. Buehler, D. Teney, M. Johnson, S. Gould, and L. Zhang, "Bottom-Up and Top-Down At-tention for Image Captioning and VQA," *arXiv:1707.07998*, 2017.

[23] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks," *arXiv:1506.01497*, Jun. 2015.

[24] J. Pennington, R. Socher, and C. D. Manning, "GloVe: Global Vectors for Word Representation," in *Empirical Methods in Natural Language Processing (EMNLP)*, 2014, pp. 1532–1543.

[25] Y. Gao, O. Beijbom, N. Zhang, and T. Darrell, "Compact bilinear pooling," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 317–326.

[26] 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 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 6904–6913.

[27] S. Antol, A. Agrawal, J. Lu, M. Mitchell, C. L. Zitnick, D. Batra, and D. Parikh, "VQA: Visual Question Answering," in *IEEE International Conference on Computer Vision (CVPR)*, 2015, pp. 2425–2433.

[28] T.-Y. Lin, M. Maire, S. Belongie, 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*, 2014, pp. 740–755.

[29] D. P. Kingma and J. Ba, "Adam: A Method for Stochastic Optimization," *arXiv:1412.6980*, 2014.

[30] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks," in *AISTATS*, 2010, pp. 249–256.

[31] N. Kalchbrenner, I. Danihelka, and A. Graves, "Grid long short-term memory," *arXiv:1507.01526*, 2015.

[32] Y. Lin, Z. Pang, D. Wang, and Y. Zhuang, "Task-driven Visual Saliency and Attention-based Visual Question Answering," *arXiv:1702.06700*, 2017.

[33] P. Wang, Q. Wu, C. Shen, and A. van den Hengel, "The VQA-Machine: Learning How to Use Existing Vision Algorithms to Answer New Questions," *arXiv:1612.05386*, 2016.

[34] I. Ilievski and J. Feng, "A Simple Loss Function for Improving the Convergence and Accuracy of Visual Question Answering Models," *arXiv:1708.00584*, 2017.

[35] Z. Yu, J. Yu, C. Xiang, J. Fan, and D. Tao, "Beyond Bilinear: Generalized Multi-modal Factorized High-order Pooling for Visual Question Answering," *arXiv:1708.03619*, 2017.

[36] D. Teney, P. Anderson, X. He, and A. van den Hengel, "Tips and Tricks for Visual Question Answering: Learnings from the 2017 Challenge," *arXiv:1708.02711*, 2017.

[37] G. Montavon, W. Samek, and K.-R. Müller, "Methods for interpreting and understanding deep neural networks," *Digital Signal Processing*, vol. 73, pp. 1–15, 2018.

[38] L. Arras, G. Montavon, K.-R. Müller, and W. Samek, "Explaining recurrent neural network predictions in sentiment analysis," in *EMNLP'17 Workshop on Computational Approaches to Subjectivity, Sentiment & Social Media Analysis (WASSA)*, 2017, pp. 159–168.