

# Faithful Multimodal Explanation for Visual Question Answering

**Jialin Wu**

Department of Computer Science  
University of Texas at Austin  
jialinwu@cs.utexas.edu

**Raymond J. Mooney**

Department of Computer Science  
University of Texas at Austin  
mooney@cs.utexas.edu

## Abstract

AI systems’ ability to explain their reasoning is critical to their utility and trustworthiness. Deep neural networks have enabled significant progress on many challenging problems such as visual question answering (VQA). However, most of them are opaque black boxes with limited explanatory capability. This paper presents a novel approach to developing a high-performing VQA system that can elucidate its answers with integrated textual and visual explanations that faithfully reflect important aspects of its underlying reasoning process while capturing the style of comprehensible human explanations. Extensive experimental evaluation demonstrates the advantages of this approach compared to competing methods using both automated metrics and human evaluation.

## 1 Introduction

Deep neural networks have made significant progress on visual question answering (VQA), the challenging AI problem of answering natural-language questions about an image (Antol et al., 2015). However successful systems (Fukui et al., 2016; Anderson et al., 2018; Yang et al., 2016; Wu et al., 2018a; Jiang et al., 2018) based on deep neural networks are difficult to comprehend because of many layers of abstraction and a large number of parameters. This makes it hard to develop user trust. Partly due to the opacity of current deep models, there has been a recent resurgence of interest in *explainable AI*, systems that can explain their reasoning to human users. In particular, there has been some recent development of explainable VQA systems (Selvaraju et al., 2017; Park et al., 2018; Hendricks et al., 2016, 2018).

One approach to explainable VQA is to generate *visual explanations*, which highlight image regions that most contributed to the system’s answer, as determined by attention mechanisms (Lu



Question: What sport is pictured? Answer: Surfing  
Explanation: Because the **man** is riding a wave on a **surfboard**.

Figure 1: Example of our multimodal explanation. It highlights relevant image regions together with a textual explanation with corresponding words in the same color.

et al., 2016) or gradient analysis (Selvaraju et al., 2017). However, such simple visualizations do not explain *how* these regions support the answer. An alternate approach is to generate a *textual explanation*, a natural-language sentence that provides reasons for the answer. Some recent work has generated textual explanations for VQA by training a recurrent neural network (RNN) to directly mimic examples of human explanations (Hendricks et al., 2016; Park et al., 2018). A *multimodal* approach that integrates *both* a visual and textual explanation provides the advantages of both. Words and phrases in the text can point to relevant regions in the image. An illustrative explanation generated by our system is shown in Figure 1.

Recent research on such multimodal VQA explanation is presented in (Park et al., 2018) that employs a form of “post hoc justification” that does not truly follow and reflect the system’s actual processing. We believe that explanations should more faithfully reflect the actual processing of the underlying system in order to provide users with a deeper understanding of the system, increasing trust for the right reasons, rather than

trying to simply convince them of the system’s reliability (Bilgic and Mooney, 2005). In order to be faithful, the textual explanation generator should focus on the set of objects that contribute to the predicted answers, and receive proper supervision from only the gold standard explanations that are consistent with the actual VQA reasoning process. Towards this end, our explanation module directly uses the VQA-attended features and is trained to only generate human explanations that can be traced back to the relevant object set using a gradient-based method called GradCAM (Selvaraju et al., 2017). To enforce local faithfulness, we also align the gradient-based visual explanations generated by the textual explanation module and the VQA module during training.

In addition, our explanations provide direct links between terms in the textual explanation and segmented items in the image, as shown in Figure 1. The result is a nice synthesis of a faithful explanation that highlights concepts actually used to compute the answer and a comprehensible, human-like, linguistic explanation. Below we describe the details of our approach and present extensive experimental results on the VQA-X (Park et al., 2018) dataset that demonstrates the advantages of our approach compared to prior work using this data (Park et al., 2018) in terms of both automated metrics and human evaluation. Further, in order to evaluate the faithfulness, we design two metrics: (1) We first measure the degree of similarity between the highlighted image segments in our textual explanations and the influential segments determined by the LIME explainer (Ribeiro et al., 2016); (2) we also measure the consistency between the gradient-based visual explanation (Selvaraju et al., 2017) of the predicted answer and the generated textual explanation.

## 2 Related Work

In this section, we review related work including visual and textual explanation generation and VQA.

### 2.1 VQA

Answering visual questions (Antol et al., 2015) has been widely investigated in both the NLP and computer vision communities. Most VQA models (Fukui et al., 2016; Lu et al., 2016) embed images using a CNN and questions using an RNN and then use these embeddings to train an answer

classifier to predict answers from a pre-extracted set. Attention mechanisms are frequently applied to recognize important visual features and filter out irrelevant parts. A recent advance is the use of the Bottom-Up-Top-Down (Up-Down) attention mechanism (Anderson et al., 2018) that attends over high-level objects instead of convolutional features to avoid emphasis on irrelevant portions of the image. We adopt this mechanism, but replace object detection (Ren et al., 2015) with segmentation (Hu et al., 2018) to obtain more precise object boundaries.

### 2.2 Visual Explanation

A number of approaches have been proposed to visually explain decisions made by vision systems by highlighting relevant image regions. GradCAM (Selvaraju et al., 2017) analyzes the gradient space to find visual regions that most affect the decision. Attention mechanisms in VQA models can also be directly used to determine highly-attended regions and generate visual explanations. Unlike conventional visual explanations, ours highlight segmented objects that are linked to words in an accompanying textual explanation, thereby focusing on more precise regions and filtering out noisy attention weights.

### 2.3 Textual and Multimodal Explanation

Visual explanations highlight key image regions behind the decision, however, they do not explain the reasoning process and crucial relationships between the highlighted regions. Therefore, there has been some work on generating textual explanations for decisions made by visual classifiers (Hendricks et al., 2016). As mentioned in the introduction, there has also been some work on multimodal explanations that link textual and visual explanations (Park et al., 2018). A recent extension of this work (Hendricks et al., 2018) first generates multiple textual explanations and then filters out those that could not be grounded in the image. We argue that a good explanation should focus on referencing visual objects that actually influenced the system’s decision, therefore generating more faithful explanations.

## 3 Approach

Our goal is to generate more faithful multimodal explanations that specifically include the segmented objects in the image that are the focus of

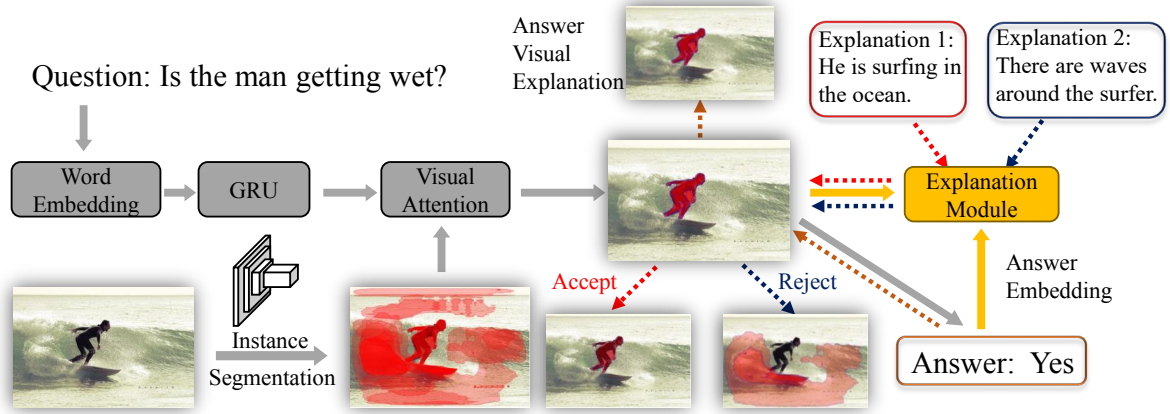


Figure 2: Model overview: We first segment the image and then predict the answer for the visual question with a pretrained VQA module. Then, we learn to embed the question, answer, and the VQA-attended features to generate textual explanations. During training, we only use the faithful human explanation whose gradient-based visual explanation is consistent with that of the predicted answer. In the example, our explanation module is only trained to generate “Explanation 1” and further enforces the consistency between this explanation and the predicted answer. “Explanation 2” is filtered out since its visual explanation is mainly focused on the waves and is not consistent with VQA module’s focus on the surfer. Dashed arrows denote gradients, gray and yellow arrows denote fixed and trainable parameters, respectively. The three smaller images denote the visual explanations for the predicted answer and the two textual explanations.

the VQA module. Figure 2 illustrates our model’s pipeline in the training phase, consisting of the VQA module (Section 3.2), and textual explanation module (Section 3.4). We first segment the objects in the image and predict the answer using the VQA module, which has an attention mechanism over those objects. Next, the explanation module is trained to generate textual explanations conditioned on the question, answer, and VQA-attended features. To faithfully train the explanation module, we filter out human textual explanations whose gradient-based visual explanation is not consistent with that of the predicted answer. For example, in Figure 2 “Explanation 1” is accepted as the textual explanation since it is mainly focused on the surfer and “Explanation 2” is rejected. For the consistent textual explanations, we also train the explanation module to align its visual explanation with the predicted answer’s to enforce local faithfulness.

### 3.1 Notation

We use  $f$  to denote the fully-connected  $fc$  layers of the neural network, and these  $fc$  layers do not share parameters. We notate the sigmoid functions as  $\sigma$ . The subscript  $i$  indexes the elements of the segmented object sets from images. Bold letters denote vectors, overlining  $\bar{\cdot}$  denotes averaging, and  $[\cdot, \cdot]$  denotes concatenation.

### 3.2 VQA Module

We base our VQA module on Up-Down (Anderson et al., 2018) with some modifications. First, we replace the two-branch gated  $\tanh$  answer classifier with single  $fc$  layers with Leaky ReLU activation (Maas et al., 2013). In order to ground the explanations in more precise visual regions, we use instance segmentation (Hu et al., 2018) to segment objects in over 3,000 categories. Specifically, we extract at most the top  $V < 80$  objects in terms of segmentation scores and concatenate each object’s  $fc6$  representation in the bounding box classification branch and  $mask.fcn[1-4]$  features in the mask generation branch to form a 2048-d vector. This results in an image feature set  $\mathbf{V}$  containing  $V$  2048-d vectors  $\mathbf{v}_i$  for each image. We encode each question as the last hidden state  $\mathbf{q}$  of a gated recurrent unit (GRU) with 512 hidden units. We learn visual attention over all the segments  $\alpha^{vqa} \in \mathbb{R}^V$ , and use the attended visual features  $\mathbf{v}_i^q$  together with the question embedding to produce a joint representation  $\mathbf{h}$ . Then the model predicts the logits  $\mathbf{s}^{vqa}$  for each answer candidate using a 2-layer  $fc$  networks, which is passed through a sigmoid function to compute the final probabilities. For the detailed network architecture, please refer to (Anderson et al., 2018). The parameters in the VQA module are fixed during the training of the explanation module.

### 3.3 Question and Answer Embedding for Explanation Generation

As suggested in (Park et al., 2018), we also encode questions and answers as input features to the explanation module. In particular, we regard the normalized answer prediction output as a multinomial distribution, and sample one answer from this distribution at each time step. We re-embed it as a one-hot vector  $\mathbf{a}_s = \text{one-hot}(\text{multinomial}(s))$ .

$$\mathbf{u}_i = \mathbf{v}_i^q \odot f(\mathbf{a}_s) \odot f(\mathbf{q}) \quad (1)$$

Next, we element-wise multiply the embedding of  $\mathbf{q}$  and  $\mathbf{a}_s$  with  $\mathbf{v}_i^q$  to compute the joint representation  $\mathbf{u}_i$ . Note that  $\mathbf{u}$  faithfully represents the focus of the VQA process, in that it is directly derived from the VQA-attended features.

### 3.4 Explanation Generation

In this section, we describe the Explanation Module depicted by the yellow box in Figure 2. The explanation module has a two-layer-LSTM architecture whose first layer produces an attention over the  $\mathbf{u}_i$ , and whose second layer learns a representation for predicting the next word using the first layer’s features.

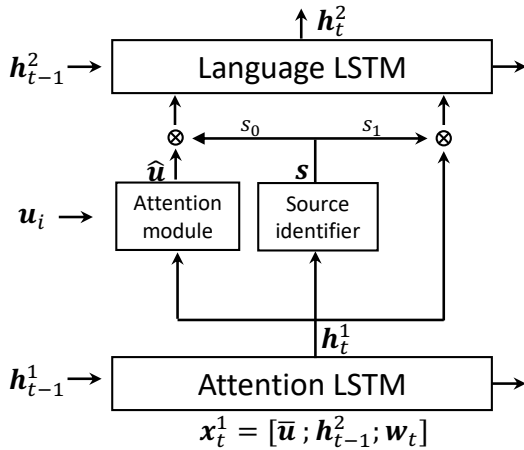


Figure 3: Overview of the explanation module.

In particular, the first visual attention LSTM takes the concatenation  $\mathbf{x}_t^1$  of the second language LSTM’s previous output  $\mathbf{h}_{t-1}^2$ , the average pooling of  $\mathbf{u}_i$ , and the previous words’ embedding as input and produces the hidden presentation  $\mathbf{h}_t^1$ . Then, an attention mechanism re-weights the image feature  $\mathbf{u}_i$  using the generated  $\mathbf{h}_t^1$  as input shown in Eq. 2. For the detailed structure, please refer to (Anderson et al., 2018).

son et al., 2018).

$$a_{i,t} = f(\text{tanh}(f(\mathbf{u}_i) + f(\mathbf{h}_t^1))) \quad (2)$$

$$\alpha_t = \text{softmax}(\alpha_t) \quad (3)$$

For the purpose of faithfully grounding the generated explanation in the image, we argue that the generator should be able to determine if the next word should be based on image content attended to by the VQA system or on learned linguistic content. To achieve this, we introduce a “source identifier” to balance the total amount of attention paid to the visual features  $\mathbf{u}_i$  and the recurrent hidden representation  $\mathbf{h}_t^1$  at each time step. In particular, given the output  $\mathbf{h}_t^1$  from the attention LSTM and the average pooling  $\bar{\mathbf{u}}_i$  over  $\mathbf{u}_i$ , we train a  $fc$  layer to produce a 2-d output  $\mathbf{s} = \sigma(f([\mathbf{h}_t^1, \bar{\mathbf{u}}_i])) = (s_0, s_1)$  that identifies which source the current generated word should be based on (*i.e.*  $s_0$  for the output of the attention LSTM<sup>1</sup> and  $s_1$  for the attended image features).

$$\mathbf{s} = \sigma(f([\mathbf{h}_t^1, \bar{\mathbf{u}}_i])) \quad (4)$$

We use the following approach to obtain training labels  $\hat{\mathbf{s}}$  for the source identifier. For each visual features  $\mathbf{u}_i$ , we assign label 1 (indicating the use of attended visual information) when there exists a segmentation  $\mathbf{u}_i$  whose cosine similarity between its category name’s GloVe representation and the current generated word’s GloVe representation is above 0.6. Given the labeled data, we train the source identifier using cross entropy loss  $\mathcal{L}_s$  as shown in Eq. 5:

$$\mathcal{L}_s = -\left(\sum_{j=0}^1 \hat{s}_j \log s_j + (1 - \hat{s}_j) \log(1 - s_j)\right) \quad (5)$$

where the  $\hat{s}_0, \hat{s}_1$  are the aforementioned labels.

Next, we concatenate the re-weighted  $\mathbf{h}_t^1$  and  $\bar{\mathbf{u}}_i$  with the output of the source identifier as the input  $\mathbf{x}_t^2 = [\mathbf{h}_t^1 s_0, \bar{\mathbf{u}}_i s_1]$  for the language LSTM. For more detail on the language LSTM structure, please refer to (Anderson et al., 2018).

With the hidden states  $\mathbf{h}_t^2$  in the Language LSTM, the output word’s probability is computed

<sup>1</sup>We tried to directly use the source weights  $s_0$  in the language LSTM’s hidden representation  $\mathbf{h}_{t-1}^2$  and found that using  $\mathbf{h}_t^1$  works better. The reason is that directly constraining  $\mathbf{h}_{t-1}^2$  makes the language LSTM forget the previously encoded content and prevents it from learning long term dependencies.

using Eq. 6:

$$p(y_t|y_{1:t-1}) = \text{softmax}(f(\mathbf{h}_t^2)) \quad (6)$$

where  $y_t$  denotes the  $t$ -th word in the explanation  $\mathbf{y}$  and  $y_{1:t-1}$  denotes the first  $t - 1$  words.

**Faithful Explanation Supervision.** Directly collecting faithful textual explanations is infeasible because it would require an annotation process where workers provide explanations based on the attended VQA features. Instead, we design an on-line algorithm that automatically filters unfaithful explanations from the human ones in the VQA-X data (Park et al., 2018) based on the idea that a proper explanation should focus on the same set of objects as the VQA module and be locally faithful. As recent research suggested that gradient-based methods more faithfully present the models’ decision making process (Zhang et al., 2018; Wu et al., 2018b; Wu and Mooney, 2019; Jain and Wallace, 2019), we define a faithfulness score  $\mathcal{S}_f$  as the cosine similarity between the Grad-CAM (Selvaraju et al., 2017) visual explanation vectors of the textual explanation and the predicted answer as shown in Eq. 7:

$$\mathcal{S}_f(\mathbf{y}) = \cos(g(s_{pred}^{vqa}, \mathbf{v}^q), g(\log p(\mathbf{y}), \mathbf{v}^q)) \quad (7)$$

where  $g$  denotes the Grad-CAM operation and the result is a vector of length  $|V|$  indicating the contribution of each segmented object.  $s_{pred}^{vqa}$  is the logit for the predicted answer.

Then, we filter out the explanations in the training set whose faithfulness scores are less than  $\xi \max(0.02 it, 1)$ , where  $\xi$  is a threshold and the  $\max(0.02 it, 1)$  term is used to jump-start the randomly initialized explanation module. For example, during training, we only accept “Explanation 1” in Figure 2 because the visual explanations of the predicted answer and the textual explanation are consistent and reject “Explanation 2”.

Since the VQA-X dataset only has explanations for the correct answers, we also discard the explanations when the predicted answers are wrong. With the remaining human explanations, we minimize the cross-entropy loss  $\mathcal{L}_{XE}$  in Eq. 8:

$$\mathcal{L}_{XE} = \sum_{t=1}^T \log(p(y_t|y_{1:t-1})) \quad (8)$$

To enforce local faithfulness, we further align these two gradient vectors using cosine distance  $\mathcal{L}_f = 1 - \mathcal{S}_f$ .

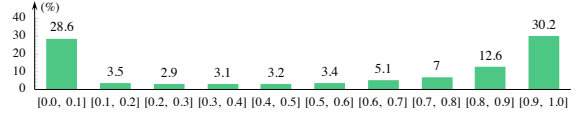


Figure 4: The distribution of explanations’ faithfulness scores in the last epoch during training.

In Figure 4, we report the distribution of the explanations’ faithfulness scores  $\mathcal{S}_f$  in the last epoch during training ( $\xi$  is set to 0.3). We observe that about 30% of the human explanations in the training set cannot be traced back to similar image segments that highly contribute to the predicted answer using our trained explanation module. These textual explanations cannot be seen as faithful either because the explanations themselves are not faithful or because the explanation module fails to develop the correct mappings between the textual explanations and the VQA-attended features. There are only a small fraction of the explanations whose faithfulness scores are in the interval of [0.1, 0.6] indicating that there is a clear boundary between whether or not an explanation is deemed faithful according to our metric.

### 3.5 Training

We pre-train the VQA module on the entire VQA v2 training set for 15 epochs using the Adam optimizer (Kingma and Ba, 2014) with a learning rate of 0.001. After that, the parameters in the VQA module are frozen. Our VQA module is capable of achieving 82.9% and 80.3% in the VQA-X train and test split respectively. and 63.5% in the VQA v2 validation set which is comparable to the baseline Up-Down model (63.2%) (Anderson et al., 2018). Note that VQA performance is not the focus of this work, and our experimental evaluation focuses on the generated explanations. Finally, we train the explanation module using the human explanations in the VQA-X dataset (Park et al., 2018) filtered for faithfulness. VQA-X contains 29,459 question answer pairs and each pair is associated with a human explanation. We train to minimize the joint loss  $\mathcal{L}$  (Eq. 9), and  $\xi$  is empirically set to 0.3. We ran the Adam optimizer for 25 epochs with a batch size of 128. The learning rate for training the explanation module is initialized to  $5e-4$  and decays by a factor of 0.8 every three epochs.

$$\mathcal{L} = \mathcal{L}_{XE} + \mathcal{L}_s + \mathcal{L}_f \quad (9)$$

				# Expl.	Textual					Visual
	$\mathcal{L}_s$	$\mathcal{F}$	$\mathcal{L}_f$		B-4	M	R-L	C	S	EMD
PJ-X (Park et al., 2018)				29K	19.5	18.2	43.7	71.3	15.1	2.64
Ours (Justification)				29K	23.5	19.0	46.2	81.2	17.2	2.46
Ours (Justification)	✓			29K	24.4	19.5	47.4	88.8	17.9	2.41
Ours (Justification)	✓			15K	24.1	18.6	46.2	83.4	16.2	2.59
Ours (Explanation)	✓	✓		15K	24.7	19.2	47.0	85.1	16.6	2.56
Ours (Explanation)	✓	✓	✓	15K	<b>25.1</b>	<b>19.7</b>	<b>48.2</b>	<b>86.7</b>	<b>17.2</b>	<b>2.52</b>

Table 1: Explanation evaluation results, the top half shows results using the entire train set and the bottom half shows results using about 15K explanations.  $\mathcal{F}$  denotes whether to filter out the unfaithful training explanations. With  $\mathcal{F}$ , the 15K explanations are the remaining explanation and without  $\mathcal{F}$ , the 15K explanations are randomly sampled from train set.  $\mathcal{L}_s$ ,  $\mathcal{L}_f$  denote the losses of the source identifier and the faithful gradient alignment, respectively. B-4, M, R-L, C and S are short hand for BLEU-4, METEOR, ROUGE-L, CIDEr and SPICE, respectively.

### 3.6 Multimodal Explanation Generation

As a last step, we link words in the generated textual explanation to image segments in order to generate the final multimodal explanation. To determine which words to link, we extract all common nouns whose source identifier weight  $s_1$  in Eq. 4 exceeds 0.5. We then link them to the segmented object with the highest attention weight  $\alpha_t$  in Eq. 2 when that corresponding output word  $y_t$  was generated, but only if this weight is greater than 0.2.<sup>2</sup>

## 4 Experimental Evaluation

This section experimentally evaluates both the textual and visual aspects of our multimodal explanations, including comparisons to competing methods and ablations that study the impact of the various components of our overall system. Finally, we present metrics and evaluation for the faithfulness of our explanations.

### 4.1 Textual Explanation Evaluation

Similar to (Park et al., 2018), we first evaluate our textual explanations using automated metrics by comparing them to the gold-standard human explanations in the VQA-X test data using standard sentence-comparison metrics: BLEU-4 (Papineni et al., 2002), METEOR (Banerjee and Lavie, 2005), ROUGE-L (Lin, 2004), CIDEr (Vedantam et al., 2015) and SPICE (Anderson et al., 2016). Table 1 reports our performance, including ablations. In particular, ‘‘Justification’’ denotes training on the entire or randomly sampled VQA-X dataset and ‘‘Explanation’’ denotes

<sup>2</sup>Due to duplicated segments, we use a lower threshold.

training only on the remaining faithful explanations. We outperform the current state-of-the-art PJ-X model (Park et al., 2018) on all automated metrics by a clear margin with only about half the explanation training data. This indicates that constructing explanations that faithfully reflect the VQA process can actually generate explanations that match human explanations better than just training to directly match human explanations, possibly by avoiding over-fitting and focusing more on important aspects of the test images.

### 4.2 Multimodal Explanation Evaluation

In this section, we present the evaluations of our model on both visual and multimodal aspects.

**Automated Evaluation:** As in previous work (Selvaraju et al., 2017; Park et al., 2018), we first used Earth Mover Distance (EMD) (Pele and Werman, 2008) to compare the image regions highlighted in our explanation to image regions highlighted by human judges. In order to fairly compare to prior results, we resize all the images in the entire test split to  $14 \times 14$  and adjust the segmentation in the images accordingly using bi-linear interpolation. Next, we sum up the multiplication of attention values and source identifiers’ values in Eq. 2 over time ( $t$ ) and assign the accumulated attention weight to each corresponding segmentation region. We then normalize attention weights over the  $14 \times 14$  resized images to sum to 1, and finally compute the EMD between the normalized attentions and the ground truth.

As shown in the Visual results in Table 1, our approach matches human attention maps more



Figure 5: Sample positively-rated explanations.

closely than PJ-X (Park et al., 2018). We attribute this improvement to the following reasons. First, our approach uses detailed image segmentation which avoids focusing on background and is much more precise than bounding-box detection. Second, our visual explanation is focused by textual explanation where the segmented visual objects must be linked to specific words in the textual explanation. Therefore, the risk of attending to unnecessary objects in the images is significantly reduced. As a result, we filter out most of the noisy attention in a purely visual explanation like that in PJ-X.

**Human Evaluation:** We also asked AMT workers to evaluate our final multimodal explanations that link words in the textual explanation directly to segments in the image. Specifically, we randomly selected 1,000 correctly answered question and asked workers “How well do the highlighted image regions support the answer to the question?” and provided them a Likert-scale set of possible answers: “Very supportive”, “Supportive”, “Neutral”, “Unsupportive” and “Completely unsupportive”. The second task was to evaluate the quality of the links between words and image regions in the explanations. We asked workers “How well do the colored image segments highlight the appropriate regions for the corresponding colored words in the explanation?” with the Like-scale choices: “Very Well”, “Well”, “Neutral”, “Not Well”, “Poorly”. We assign five questions in each AMT HIT with one “validation” item to control the HIT’s qualities.

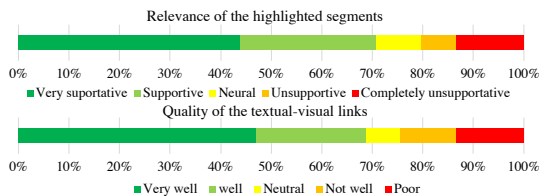


Figure 6: Human evaluation results.

As shown in Figure 6, in both cases, about 70% of the evaluations are positive and about 45% of them are strongly positive. This indicates that our multimodal explanations provide good connections among visual explanations, textual explanations, and the VQA process. Figure 5 presents some sample positively-rated multimodal explanations.

### 4.3 Faithfulness Evaluation

In this section, we measure the faithfulness of our explanations, i.e. how well they reflect the underlying VQA system’s reasoning. First, we measured how many words in a generated explanation are actually linked to a visual segmentation in the image. We analyzed the explanations from 1,000 correctly answered questions from the test data. On average, our model is able to link 1.6 words in an explanation to an image segment, indicating that the textual explanation is actually grounded in objects detected by our VQA system.

**Faithfulness Evaluation using LIME.** We use the model-agnostic explainer LIME (Ribeiro et al., 2016) to determine the segmented objects that most influenced a particular answer, and measure how well the objects referenced in our explanation match these influential segments. We regard all the detected visual segments as the “interpretable” units used by LIME to explain decisions. Using these interpretable units, LIME applies LASSO with the regularization path (Efron et al., 2004) to learn a linear model of the local decision boundary around the example to be explained. In particular, we collect 256 points around the example by randomly blinding each segment’s features with a probability of 0.4. The highly weighted features in this model are claimed to provide a faithful explanation of the decision on this example (Ribeiro et al., 2016). The complexity of the explanation is controlled by the number of units,  $K$ , that can be used in this linear model. Using the coefficients  $\mathbf{w}$  of LIME’s weighted linear model, we compare the object segments selected by LIME to the set of objects that are actually linked to words in our explanations. Specifically, we define our faithfulness metric as:

$$score = \frac{\sum_{i=1}^{|V|} |w_i| \max_{j \in \mathcal{L}} \cos(\mathbf{v}_i, \mathbf{v}_j)}{\sum_{i=1}^{|V|} |w_i|} \quad (10)$$

where  $\mathbf{v}_i$  denotes the visual feature of the  $i$ -th segmented object and the  $\mathcal{L}$  denotes the set of

explanation-linked objects. For each object in the LIME explanation, it finds the closest object in our explanation and multiplies its LIME weight by this similarity. The normalized sum of these matches is used to measure the similarity of the two explanations.

We collect all correctly answered questions in the VQA-X test set, and Table 2 reports the average score for their explanations using models trained on 15K training explanations with different numbers of interpretable units  $K$ . The influential objects recognized by LIME match objects that are linked to words in our explanations with an average cosine similarity around 0.7. This indicates that the explanations are faithfully making reference to visual segmentations that actually influenced the decision of the underlying VQA system. Also, we observe that training with faithful human explanation outperforms purely mimicking human explanations in terms of our faithful metric, and further enforcing the local faithfulness achieves a better result.

	K = 1	K = 2	K = 3
Ours (Random)	0.588	0.601	0.574
Ours (Filtered)	0.636	0.651	0.643
Ours (Filtered+ $\mathcal{L}_f$ )	<b>0.686</b>	<b>0.705</b>	<b>0.678</b>

Table 2: Evaluation of LIME-based faithfulness scores for different numbers of interpretable units  $K$  using 15K training explanations. “Random” means the training explanations are randomly sampled from the train set, and “Filtered” means the models are trained using the remaining faithful explanations.

**Faithfulness Evaluation using Grad-CAM.** We also evaluated the consistency between the Grad-CAM visual explanation vectors from the textual explanation and the predicted answer using the faithful score  $\mathcal{S}_f$  defined in Eq. 7. Table 3 reports the results from using filtered verses randomly sampled explanations for training. We observe that with faithful human explanations, the average faithfulness evaluation score increases 7% over training with randomly sampled explanations. Moreover, with the faithfulness loss  $\mathcal{L}_f$ , the model can better align the visual explanation for the textual explanation with that for the predicted answer, leading to a further 11% increase.

We also report the distribution of the generated explanations’ cosine similarity between their visual explanation and the visual explanation of the

	# Expl.	Average $\mathcal{S}_f$
Ours (Random)	15K	0.38
Ours (Filtered)	15K	0.45
Ours (Filtered+ $\mathcal{L}_f$ )	15K	<b>0.56</b>

Table 3: Average faithfulness evaluation score using different explanations models. “Random” means the training explanations are randomly sampled from the train set, and “Filtered” means the models are trained using the remaining faithful explanations.

answers in Figure 7. The fraction of the faithfulness scores between the interval  $[0.0, 0.1]$  is significantly decreased by over 17% when using the faithful human explanations for supervision and further enforcing the local faithfulness with the faithfulness loss,  $\mathcal{L}_f$ .

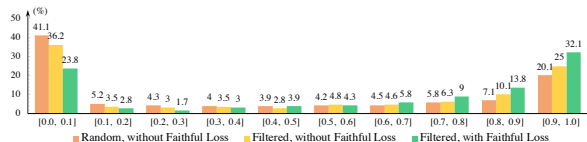


Figure 7: The distribution of explanations’ cosine similarity between the visual explanation of the textual explanation and the predicted answer.

## 5 Conclusion and Future Work

This paper has presented a new approach to generating multimodal explanations for visual question answering systems that aims to more faithfully represent the reasoning of the underlying VQA system while maintaining the style of human explanations. The approach generates textual explanations with words linked to relevant image regions actually attended to by the underlying VQA system. Experimental evaluations of both the textual and visual aspects of the explanations using both automated metrics and crowdsourced human judgments were presented that demonstrate the advantages of this approach compared to a previously-published competing method. In the future, we would like to incorporate more information from the VQA networks into the explanations. In particular, we would like to integrate *network dissection* (Bau et al., 2017) to allow recognizable concepts in the learned hidden-layer representations to be included in explanations.



## Acknowledgement

This research was supported by the DARPA XAI program under a grant from AFRL.

## References

- Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. 2016. Spice: Semantic Propositional image caption evaluation. In *European Conference on Computer Vision*, pages 382–398.
- Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. 2018. Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, volume 3, page 6.
- Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. 2015. VQA: Visual question answering. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2425–2433.
- Satanjeev Banerjee and Alon Lavie. 2005. Meteor: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments. In *Proceedings of the ACL workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pages 65–72.
- David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. 2017. Network Dissection: Quantifying Interpretability of Deep Visual Representations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Mustafa Bilgic and Raymond Mooney. 2005. [Explaining Recommendations: Satisfaction vs. Promotion](#). In *Proceedings of Beyond Personalization 2005: A Workshop on the Next Stage of Recommender Systems Research at the 2005 International Conference on Intelligent User Interfaces*.
- Bradley Efron, Trevor Hastie, Iain Johnstone, Robert Tibshirani, et al. 2004. Least angle regression. *The Annals of statistics*, 32(2):407–499.
- Akira Fukui, Dong Huk Park, Daylen Yang, Anna Rohrbach, Trevor Darrell, and Marcus Rohrbach. 2016. Multimodal Compact Bilinear Pooling for Visual Question Answering and Visual Grounding. In *Proceedings of the 2016 Empirical Methods on Natural Language Processing*.
- Lisa Anne Hendricks, Zeynep Akata, Marcus Rohrbach, Jeff Donahue, Bernt Schiele, and Trevor Darrell. 2016. Generating Visual Explanations. In *European Conference on Computer Vision*, pages 3–19. Springer.
- Lisa Anne Hendricks, Ronghang Hu, Trevor Darrell, and Zeynep Akata. 2018. Grounding Visual Explanations. *arXiv preprint arXiv:1807.09685*.
- Ronghang Hu, Piotr Dollár, Kaiming He, Trevor Darrell, and Ross Girshick. 2018. Learning to Segment Every Thing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Sarthak Jain and Byron C Wallace. 2019. Attention is not explanation. *arXiv preprint arXiv:1902.10186*.
- Yu Jiang, Vivek Natarajan, Xinlei Chen, Marcus Rohrbach, Dhruv Batra, and Devi Parikh. 2018. Pythia v0. 1: the winning entry to the vqa challenge 2018. *arXiv preprint arXiv:1807.09956*.
- Diederik P Kingma and Jimmy Ba. 2014. Adam: A Method for Stochastic Optimization. *Proceedings of the 3rd International Conference on Learning Representations*.
- Chin-Yew Lin. 2004. Rouge: A Package for Automatic Evaluation of Summaries. *Text Summarization Branches Out*.
- Jiasen Lu, Jianwei Yang, Dhruv Batra, and Devi Parikh. 2016. Hierarchical Question-Image Coattention for Visual Question Answering. In *Advances In Neural Information Processing Systems*, pages 289–297.
- Andrew L Maas, Awni Y Hannun, and Andrew Y Ng. 2013. Rectifier Nonlinearities Improve Neural Network Acoustic Models. In *International Conference on Machine Learning*, volume 30, page 3.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. [Bleu: A Method for Automatic Evaluation of Machine Translation](#). In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics, ACL '02*, pages 311–318, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Dong Huk Park, Lisa Anne Hendricks, Zeynep Akata, Anna Rohrbach, Bernt Schiele, Trevor Darrell, and Marcus Rohrbach. 2018. Multimodal Explanations: Justifying Decisions and Pointing to the Evidence. In *Proceedings of the IEEE conference on computer vision and pattern recognition*.
- Ofir Pele and Michael Werman. 2008. A Linear Time Histogram Metric for Improved Sift Matching. In *Computer Vision—ECCV 2008*, pages 495–508. Springer.
- Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2015. Faster R-CNN: Towards Real-time Object Detection with Region Proposal Networks. pages 91–99.
- Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. Why should i trust you?: Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144. ACM.

- Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. 2017. Grad-cam: Visual Explanations from Deep Networks via Gradient-Based Localization. In *Proceedings of the IEEE International Conference on Computer Vision*.
- Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. 2015. Cider: Consensus-Based Image Description Evaluation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4566–4575.
- Jialin Wu, Zeyuan Hu, and Raymond J Mooney. 2018a. Joint image captioning and question answering. *arXiv preprint arXiv:1805.08389*.
- Jialin Wu, Dai Li, Yu Yang, Chandrajit Bajaj, and Xi-angyang Ji. 2018b. Dynamic Filtering with Large Sampling Field for Convnets. *ECCV*.
- Jialin Wu and Raymond J Mooney. 2019. Self-critical reasoning for robust visual question answering. *arXiv preprint arXiv:1905.09998*.
- Zichao Yang, Xiaodong He, Jianfeng Gao, Li Deng, and Alex Smola. 2016. Stacked attention networks for image question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 21–29.
- Jianming Zhang, Sarah Adel Bargal, Zhe Lin, Jonathan Brandt, Xiaohui Shen, and Stan Sclaroff. 2018. Top-down neural attention by excitation back-prop. *International Journal of Computer Vision*, 126(10):1084–1102.