

# Cycle-Consistency for Robust Visual Question Answering

Meet Shah<sup>1</sup>, Xinlei Chen<sup>1</sup>, Marcus Rohrbach<sup>1</sup>, Devi Parikh<sup>1,2</sup>

<sup>1</sup>Facebook AI Research, <sup>2</sup>Georgia Institute of Technology

{meetshah, xinleic, mrf}@fb.com, dparikh@gatech.edu

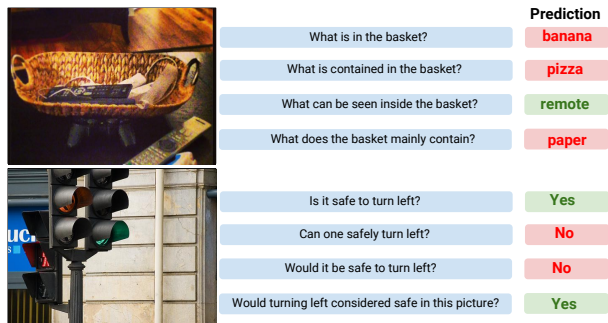
## Abstract

Despite significant progress in Visual Question Answering over the years, robustness of today’s VQA models leave much to be desired. We introduce a new evaluation protocol and associated dataset (VQA-Rephrasings) and show that state-of-the-art VQA models are notoriously brittle to linguistic variations in questions. VQA-Rephrasings contains 3 human-provided rephrasings for 40k questions spanning 40k images from the VQA v2.0 validation dataset. As a step towards improving robustness of VQA models, we propose a model-agnostic framework that exploits cycle consistency. Specifically, we train a model to not only answer a question, but also generate a question conditioned on the answer, such that the answer predicted for the generated question is the same as the ground truth answer to the original question. Without the use of additional annotations, we show that our approach is significantly more robust to linguistic variations than state-of-the-art VQA models, when evaluated on the VQA-Rephrasings dataset. In addition, our approach outperforms state-of-the-art approaches on the standard VQA and Visual Question Generation tasks on the challenging VQA v2.0 dataset.

## 1. Introduction

Visual Question Answering (VQA) applications allow a human user to ask a machine questions about images – be it a user interacting with a visual chat-bot or a visually impaired user relying on an assistive device. As this technology steps out of the realm of curated datasets towards real-world settings, it is desirable that VQA models be robust to and consistent across reasonable variations in the input modalities. While there has been significant progress in VQA over the years [1, 18, 2, 10, 21, 46, 4, 5], today’s VQA models are, however, far from being robust.

VQA is a task that lies at the intersection of language and vision. Existing works have studied the robustness and sensitiveness of VQA models to meaningful semantic variations in images [10], changing answer distributions [2] and adversarial attacks [44] to images. However, to the best of our knowledge, no work has studied the robustness of VQA



	Prediction
What is in the basket?	banana
What is contained in the basket?	pizza
What can be seen inside the basket?	remote
What does the basket mainly contain?	paper
Is it safe to turn left?	Yes
Can one safely turn left?	No
Would it be safe to turn left?	No
Would turning left considered safe in this picture?	Yes

Figure 1. **Existing VQA models are brittle.** Shown above are examples from our new large-scale **VQA-Rephrasings** dataset that enables systematic evaluation of robustness of VQA models to linguistic variations in the input question. Also shown are answers predicted by a state-of-the-art VQA model [46]. We see that the model predicts different answers for different reasonable rephrasings of the same question. We propose a novel model-agnostic framework that exploits cycle consistency in question answering and question generation to make VQA models more robust, without using additional annotation. Moreover, it outperforms state-of-the-art models on the standard VQA and Visual Question Generation tasks on the VQA v2.0 dataset.

models to linguistic variations in the input question. This is important both from the perspective of VQA being a benchmark to test multi-modal AI capabilities (do our VQA models really “understand” the question when answering it?) and for applications (human users are likely to phrase the same query in a variety of different linguistic forms). However, today’s state-of-the-art VQA models are brittle to such linguistic variations as can be seen in Fig. 1.

One approach to make VQA models more robust is to collect a dataset with diverse rephrasings of questions to train VQA models. This requires additional human annotation and thus is not always scalable in real-world settings. Alternatively, an automatic approach that does not require additional human intervention but results in a VQA model that is more robust to linguistic variations observed in the natural language open-ended questions is desirable.

We propose a novel model-agnostic framework that re-

lies on cycle consistency to learn robust VQA models without requiring additional annotation. Specifically, we train the model to not just answer a question, but also to generate diverse, semantically similar variations of questions conditioned on the answer. We enforce that the answer predicted for a generated question matches the ground truth answer to the original question. In other words, the model is being trained to predict the same (correct) answer for a question and its (generated) rephrasing.

Advantages of our proposed approach are two fold. First, enforcing consistent correctness across diverse rephrasings allows models to generalize to unseen semantically equivalent variations of questions at test time. The model achieves this by generating linguistically diverse rephrasings of questions on-the-fly and training with these variations. Second, a model trained generatively to generate a valid question given a candidate answer and image has a stronger multi-modal understanding of vision and language. Questions tend to have less learnable biases [28]. As a result, models that can jointly perform the task of question generation and question answering are less prone to taking “shortcuts” and exploiting linguistic priors in questions. Indeed, we find that models trained with our approach outperform existing state-of-the-art models on both VQA and Visual Question Generation (VQG) tasks on VQA v2.0 [10].

We also observed that one reason for limited development of VQA models robust to linguistic variations in input questions is due to the lack of a benchmark to measure robustness. A lack of such a benchmark makes it hard to quantitatively realize the inflated capabilities and limited multi-modal understanding of modern VQA models and consequently inhibits progress in pushing the state-of-the-art in multi-modal understanding aspects of computer vision. To enable quantitative evaluation of robustness and consistency of VQA models across linguistic variations in input questions, we collect a large-scale dataset – **VQA-Rephrasings** (Section 4) based of the VQA v2.0 dataset [10]. VQA-Rephrasings contains 3 human-provided rephrasings for  $\sim 40k$  questions on  $\sim 40k$  images from the validation split of the VQA v2.0 dataset. We also propose metrics to measure the robustness of VQA models across different question rephrasings. Further, we benchmark several state-of-the-art VQA models [4, 6, 21, 46] on our proposed VQA-Rephrasings dataset to highlight the fragility of VQA models to question rephrasings. We observe a significant drop when VQA models are required to be consistent in addition to being correct (Section 5), which reinforces our belief that existing VQA models do not understand language “enough”. We show that VQA models trained with our approach are significantly more robust across question rephrasings than their existing counterparts on the proposed VQA-Rephrasings dataset.

In this paper, our contributions are the following:

- We propose a model-agnostic cycle-consistent training scheme that enables VQA models to be more robust to linguistic variations observed in natural language open-ended questions.
- To evaluate the robustness of VQA models to linguistic variations, we introduce a large-scale **VQA-Rephrasings** dataset and an associated consensus score. VQA-Rephrasings consists of 3 rephrasings for  $\sim 40k$  questions on  $\sim 40k$  images from the VQA v2.0 validation dataset, resulting in a total of  $\sim 120k$  questions rephrasing by humans.
- We show that models trained with our approach outperform state-of-the-art on the standard VQA and Visual Question Generation tasks on the VQA v2.0 dataset and are significantly more robust to linguistic variations on VQA-Rephrasings.

## 2. Related Work

**Visual Question Answering.** There has been tremendous progress in building models for VQA using LSTMs [14] and convolutional networks [24]. VQA models spanning different paradigms like attention networks [45, 21], module networks [15, 5, 18], relational networks [35] and multi-modal fusion [6] have been proposed. Our method is model-agnostic and is applicable with any existing VQA architecture.

**Robustness.** Robustness of VQA models has been studied in several contexts [2, 44, 10]. For example, [2] studies the robustness of VQA models to changes in the answer distributions across training and test settings; [47] analyzes the extent of visual grounding in VQA models by studying robustness of VQA models to meaningful semantic changes in images; [44] shows that despite the use of an advanced attention mechanism, it is easy to fool a VQA model with very minor changes in the image. Our work, however, aims to complete the study in robustness by benchmarking and improving robustness of VQA models to linguistic and compositional variations in questions in the form of rephrasings. Robustness has also been studied in natural language processing (NLP) systems [8, 13] in contexts of bias [39, 38], domain-shift [25] and syntactic variations [16]. To counter these issues in NLP systems, solutions like linguistically motivated data-augmentation [25] and adversarial training [16] have been proposed. We study this in the context of visual question answering which is a multi-modal task which grounds language into the visual world.

**(Visual) Question Generation.** Question Generation (QG) as a task has been studied extensively by [3, 20, 37, 43] in NLP. Generating questions conditioned on an image was introduced in [32] and a large-scale VQG dataset was collected by [33] to evaluate visually grounded question

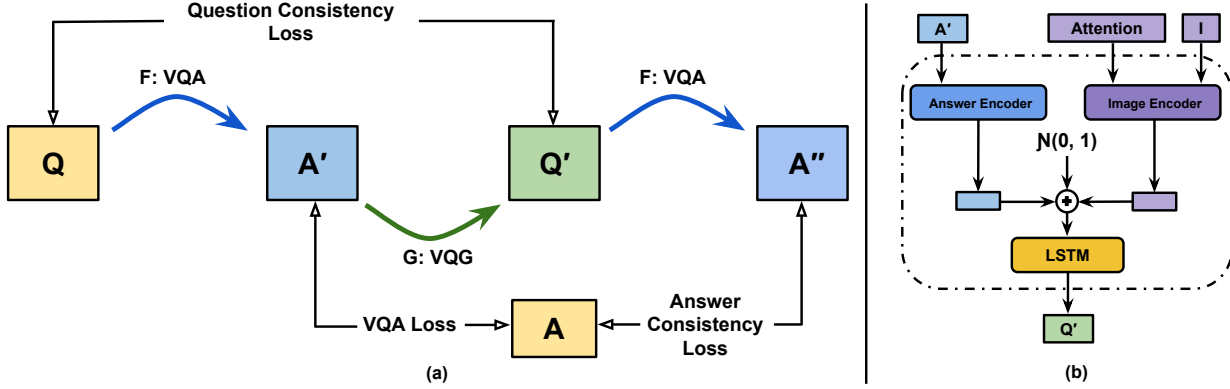


Figure 2. (a) **Abstract representation of the proposed cycle-consistent training scheme:** Given a triplet of image  $I$ , question  $Q$ , and ground truth answer  $A$ , a VQA model is a transformation  $F : (Q, I) \mapsto A'$  used to predict the answer  $A'$ . Similarly, a VQG model  $G : (A', I) \mapsto Q'$  is used to generate a rephrasing  $Q'$  of  $Q$ . The generated rephrasing  $Q'$  is passed through  $F$  to obtain  $A''$  and consistency is enforced between  $Q$  and  $Q'$  and between  $A'$  and  $A''$ . Image  $I$  is not shown for clarity. (b) **Detailed architecture of our visual question generation module  $G$ .** The predicted answer  $A'$  and image  $I$  are embedded to a lower dimension using task-specific encoders and the resulting feature maps are summed up with additive noise and fed to an LSTM to generate questions rephrasings  $Q'$ .

generation capabilities of models. More recently, there has been work on generating questions that are diverse [17, 45]. Training models to ask informative questions about an image in an active learning fixed-budget setting was explored in [31]. While these techniques generate questions about an image in an answer-agnostic manner, techniques like [28] propose a variational LSTM based model trained with reinforcement learning to generate answer-specific questions for an image. More recently, [26] generates answer-specific questions for specific question-types by modelling question generation as a dual task of question answering. Unlike [26], our method is not restricted to generating questions only for specific question types. Different from previous works, the goal of our VQG component is to automatically generate question rephrasings that make the VQA models more robust to linguistic variations. To the best of our knowledge, we are the first to demonstrate that the VQG module can be used to improve VQA accuracy in a cycle-consistent setting.

**Cycle-Consistent Learning.** Using cycle-consistency to regularize the training of models has been used extensively in object tracking [40], machine translation [11], unpaired image-to-image translation [48] and text-based question answering [41]. Consistency enables learning of robust models by regularizing transformations that map one interconnected modality or domain to the other. While cycle consistency has been used vastly in the domains involving a single modality (text-only or image-only), it hasn't been explored in the context of multi-modal tasks like VQA. Cycle-consistency in VQA can be also thought of as an online data-augmentation technique where the model is trained on several generated rephrasings of the same question.

### 3. Approach

We now introduce our cycle-consistent scheme to train robust VQA models. Given a triplet of image  $I$ , question  $Q$ , and ground truth answer  $A$ , a generic VQA model can be formulated as a transformation  $F : (Q, I) \mapsto A'$ , where  $A'$  is the answer predicted by the model as in Fig. 2(a). Similarly, a generic VQG model can be formulated as a transformation  $G : (A', I) \mapsto Q'$  as in Fig. 2(b). For a given  $(I, Q, A)$  triplet, we first obtain an answer prediction  $A'$  using the VQA model  $F$  for the original question  $Q$ . We then use the predicted answer  $A'$  and the image  $I$  to generate a question  $Q'$  which is semantically similar to  $Q$  using the VQG model  $G$ . Lastly, we obtain a answer prediction  $A''$  for the generated question  $Q'$ .

Our design of consistency components is inspired by two beliefs. Firstly, a model which can generate a semantically and syntactically correct question given an answer and an image, has a better understanding of the cross-modal connections among the image, the question and the answer, which make them a valid  $(I, Q, A)$  triplet. Secondly, assuming the generated question  $Q'$  is a valid rephrasing of the original question, a robust VQA model should answer this rephrasing with the same answer as the original question  $Q$ . In practice, however, there are several challenges that inhibit enforcement of cycle-consistency in VQA. We discuss these challenges and describe the key components of our framework geared to tackle them in the following sections.

#### 3.1. Question Generation Module

Since VQA is a setting where there is high disparity in the information content of involved modalities (a question and answer pair is a very lossy compressed representation of

the image), learning transformations that map one modality to another is non-trivial. In cycle-consistent models dealing with single-modalities, transformations need to be learned across different domains of the same modality (image or text) with roughly similar information contents. However in a multi-modality transformation like VQG, learning a transformation from a low information modality (such as answer) to high information modality (question) needs additional supervision. We provide this additional supervision to the VQG model in the form of attention. To generate a rephrasing  $Q'$ , the VQG is guided to attend at regions of the image which were used by the VQA model to answer the original question  $Q$ . Unlike [26], this enables our models to generate questions more similar to the original question from answers like “yes”, which could possibly have a large space of plausible questions.

We model the question generation module  $G$  in a fashion similar to a conditional image captioning model. The question generation module consists of two linear encoders that transform attended image features obtained from VQA model and the distribution over answer space to lower dimensional feature vectors. We sum these feature vectors with additive noise and pass them through an LSTM which is trained to reconstruct the original question and optimized by minimizing the negative log likelihood with teacher-forcing. Note that unlike [28, 26] we do not pass the one-hot vector representing the answer obtained, or an embedding of the answer obtained to the question generation, but rather the predicted distribution over answers. This enables the question generation module to learn to map the model’s confidence over answers to the generated question.

Throughout the paper, **Q-consistency** implies addition of a VQG module  $G$  on top of the base VQA model  $F$  to generate rephrasings  $Q'$  from the image  $I$  and the predicted answer  $A'$  with an associated Q-consistency loss  $\mathcal{L}_G(Q, Q')$ . Similarly, **A-consistency** implies passing all questions generated  $Q'$  by the VQG Model  $G$  to the VQA model  $F$  and an associated A-consistency loss  $\mathcal{L}_{cycle}(A, A'')$ . The overall loss can be written as:

$$\mathcal{L}_{total} = \mathcal{L}_F(A, A') + \lambda_G \mathcal{L}_G(Q, Q') + \lambda_C \mathcal{L}_{cycle}(A, A'') \quad (1)$$

where  $\mathcal{L}_F(A, A')$  and  $\mathcal{L}_{cycle}(A, A'')$  (*i.e.* A-Consistency Loss) are cross-entropy losses,  $\mathcal{L}_G(Q, Q')$  (*i.e.* Q-Consistency Loss) is sequence generation loss [30] and  $\lambda_G, \lambda_C$  are tunable hyperparameters.

### 3.2. Gating Mechanism

One of the assumptions of our proposed cycle-consistent training scheme is that the generated question is always semantically and syntactically correct. However, in practice this is not always true. Previous attempts [19] at naively generating questions conditioned on the answer and using

them without filtering to augment the training data have been unsuccessful. Like the visual question answering module, the visual question generation module is also not perfect. Therefore not all questions generated by the question generator are coherent and consistent with the image, the answer and the original question. To overcome this issue, we propose a gating mechanism, which automatically filters undesirable questions generated by the VQG model before passing them to the VQA model for A-consistency. The gating mechanism is only relevant when used in conjunction with A-consistency. We retain only those questions which either the VQA model  $F$  can answer correctly or have a cosine similarity with the original question encoding greater than a threshold  $T_{sim}$ .

### 3.3. Late Activation

One key component of designing cycle consistent models is to prevent mode collapse. Learning cycle-consistent models in complex settings like VQA needs a carefully chosen training scheme. Since cycle-consistent models have several interconnected sub-networks learning different transformations, it is important to ensure that each of these sub-networks are working in harmony. For example, if the VQA model  $F$  and VQG model  $G$  are jointly trained and consistency is enforced in early stages of training, it is possible that both models can just “cheat” by both producing undesirable outputs. We overcome this by activating cycle-consistency at later stages of training, to make sure both VQA and VQG models have been sufficiently trained to produce reasonable outputs. Specifically, we enable the loss associated with cycle-consistency after a fixed  $A_{iter}$  iterations in the training process.

We find these design choices for question generation module, gating mechanism and late activation to be crucial for effectively training our model. We demonstrate this empirically via ablation studies in Table 2. As we want to increase the robustness of the VQA model to all generated variations, the weights between VQA models which answer the original question and the generated rephrasing are shared. Our formulation of cycle-consistency in VQA can be also thought of as an online data-augmentation technique where the model is trained on several generated rephrasings of the same question and hence is more robust to such anomalies during inference. We show that with clever training strategy, coupled with attention and carefully chosen model architectures for question generation, incorporating cycle consistency for VQA is possible and not only leads to models that are better performing, but also more robust and consistent. In addition, we show that this robustness also imparts VQA models the ability to better predict their own failures.

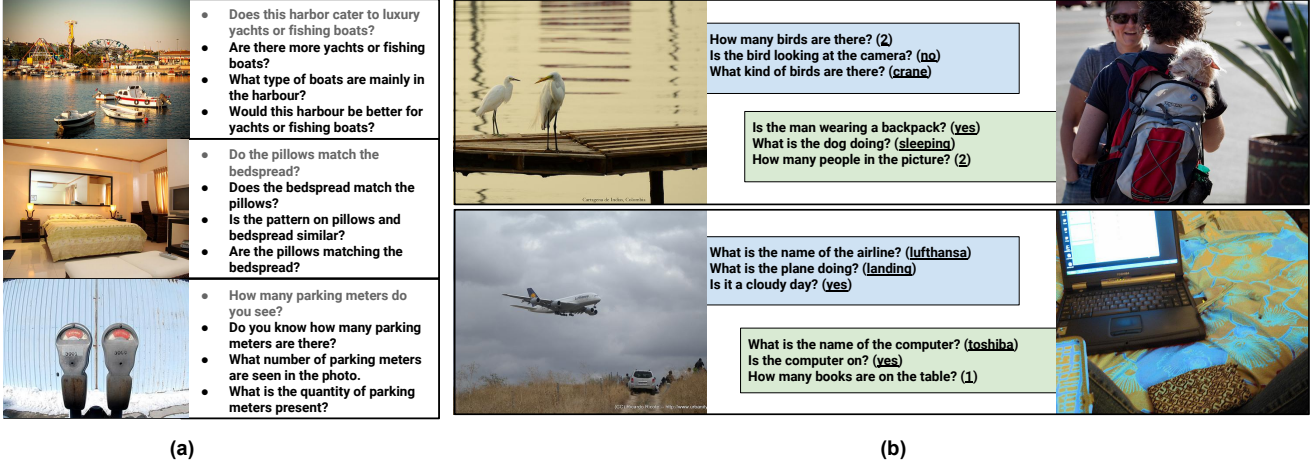


Figure 3. (a) Qualitative examples from our VQA-Rephrasings dataset. The first question (shown in gray) in each block is the original question from VQA v2.0 validation set, the questions that follow (shown in black) are rephrasings collected in VQA-Rephrasings. (b) Qualitative examples of answer conditioned question generation (input answer) by our VQG module

#### 4. VQA-Rephrasings Dataset

In this section, we introduce the VQA-Rephrasings dataset, which is the first dataset that enables evaluation of VQA models for robustness and consistency to different rephrasings of questions with the same meaning.

We use the validation split of VQA v2.0 [10] as our base dataset which contains a total of 214,354 questions spanning over 40,504 images. We randomly sample 40,504 questions (one question per image) from the base dataset to form a sampled subset. We collect 3 rephrasings of each question in the sampled subset using human annotators in two stages. In the first stage, humans were primed with the original question and the corresponding true answer and asked to rephrase the question such that answer to the rephrased question remains the same as the original answer. To ensure rephrasings from first stage are *syntactically* correct and *semantically* inline with the original question, we filter the collected responses in the next stage.

In the second stage, humans were primed with the original question and its rephrasing and were asked to label the rephrasing invalid if: (a) the plausible answer to the original question and its rephrasing is different (*i.e.* if the question and its rephrasing have different intents) or (b) if the rephrasing is grammatically incorrect. We collected 121,512 rephrasings from the original 40,504 questions in the first stage. Of these, 1320 rephrasings were flagged as invalid in the second stage and were rephrased again in the first stage. Humans were shown examples of incorrect rephrasings in the first stage to minimize the number of invalid rephrasings.

The final dataset consists of 162,016 questions (including the original 40,504 questions) spanning 40,504 images

with an average of  $\sim 3$  rephrasings per original question. A few qualitative examples from the collected dataset can be seen in Fig. 3(a). Additional details about the data collection, interfaces used and exhaustive dataset statistics can be found in Appendix A.

**Consensus Score.** Intuitively, for a VQA model to be consistent across various rephrasings of the same question, the answer to all rephrasings should be the same. We measure this by a Consensus Score  $CS(k)$ . For every group  $Q$  consisting of  $n$  rephrasings, we sample all subsets of size  $k$ . The consensus score  $CS(k)$  is defined as the ratio of the number of subsets where *all* the answers are correct and the total number of subsets of size  $k$ . The answer to a question is considered correct if it has a non-zero VQA Accuracy  $\theta$  as defined in [1].  $CS(k)$  is formally defined as:

$$CS(k) = \sum_{Q' \subset Q, |Q'|=k} \frac{S(Q')}{{}^nC_k} \quad (2)$$

$$S(Q') = \begin{cases} 1 & \text{if } \forall q \in Q' \theta(q) > 0, \\ 0 & \text{otherwise.} \end{cases} \quad (3)$$

Where  ${}^nC_k$  is number of subsets of size  $k$  sampled from a set of size  $n$ . As consensus score is a all-or-nothing score, to achieve a non-zero consensus score at  $k$  for a group of questions  $Q$ , the model has to answer at least  $k$  questions correctly in a group of questions  $Q$ . When  $k = |Q|$  (*e.g.* when  $k = 4$  in VQA-Rephrasings), the model needs to answer all rephrasings of a question and the original question correctly in order to get a non-zero consensus score. It is evident that a model with higher average consensus score at high values of  $k$  is quantitatively more robust to linguistic variations in questions than a model with a lower score.



Model	CS(k)				VQA Accuracy	
	k=1	k=2	k=3	k=4	ORI	REP
MUTAN [6]	56.68	43.63	38.94	32.76	59.08	46.87
BUTD [4]	60.55	46.96	40.54	34.47	61.51	51.22
BUTD + CC	<b>61.66</b>	<b>50.79</b>	<b>44.68</b>	<b>42.55</b>	<b>62.44</b>	<b>52.58</b>
Pythia [46]	63.43	52.03	45.94	39.49	64.08	54.20
Pythia + CC	<b>64.36</b>	<b>55.45</b>	<b>50.92</b>	<b>44.30</b>	<b>64.52</b>	<b>55.65</b>
BAN [21]	64.88	53.08	47.45	39.87	64.97	55.87
BAN + CC	<b>65.77</b>	<b>56.94</b>	<b>51.76</b>	<b>48.18</b>	<b>65.87</b>	<b>56.59</b>

Table 1. **Consensus performance on VQA-Rephrasings dataset.** CS(k) as defined in Eq. 2 is consensus score which is non-zero only if *at least*  $k$  rephrasings are answered correctly, zero otherwise; averaged across all group of questions. ORI represent a split of questions from VQA-Rephrasings which are original questions from VQA v2.0 and their corresponding rephrasings are represented by the split REP. Models trained with our cycle-consistent (CC) framework consistently outperform their baseline counterparts at all values of  $k$ .

## 5. Experiments

### 5.1. Consistency Performance

We start by benchmarking a variety of existing VQA models on our proposed VQA-Rephrasings dataset.

**MUTAN** [6] <sup>1</sup> parametrizes bilinear interactions between visual and textual representations using a multi-modal low-rank decomposition. MUTAN uses skip-thought [22] sentence embeddings to encode the question and Resnet-152 [12] to encode images. MUTAN achieves 63.20% accuracy on VQA v2.0 test-dev. Among all models we analyze, MUTAN is the only model which uses sentence embeddings to encode questions and Resnet to encode images.

**Bottom-Up Top-Down Attention (BUTD)** [4] <sup>2</sup> incorporates bottom-up attention in VQA by extracting features associated with image regions proposed by Faster-RCNN [36] pretrained on Visual Genome [23]. BUTD model won the VQA Challenge in 2017 and achieves 66.25% accuracy on VQA v2.0 test-dev.

**Pythia** [46] <sup>3</sup> extends the BUTD model by incorporating co-attention [29] between question and image regions. Pythia uses features extracted from Detectron [9] pretrained on Visual Genome. An ensemble of Pythia models won the VQA Challenge in 2018 using additional training data from Visual Genome [23] and using additional Resnet[12] features. In this study, we use Pythia models which do not use Resnet features. Pythia without using Resnet features,

<sup>1</sup><https://github.com/Cadene/vqa.pytorch>

<sup>2</sup><https://github.com/hengyuan-hu/bottom-up-attention-vqa>

<sup>3</sup><https://github.com/facebookresearch/pythia>

Model	val	test-dev
MUTAN [6]	61.04	63.20
BUTD [4]	65.05	66.25
+ Q-consistency	65.38	66.83
+ A-consistency	60.84	62.18
+ Gating	<b>65.53</b>	<b>67.55</b>
Pythia [46]	65.78	68.43
+ Q-consistency	65.39	68.58
+ A-consistency	62.08	63.77
+ Gating	<b>66.03</b>	<b>68.88</b>
BAN [21]	66.04	69.64
+ Q-consistency	66.27	69.69
+ A-consistency	64.96	66.31
+ Gating	<b>66.77</b>	<b>69.87</b>

Table 2. **VQA Performance and ablation studies on VQA v2.0 validation and test-dev splits.** Each row in blocks represents a component of our cycle-consistent framework added to the previous row. First row in each block represents the baseline VQA model  $F$ . Q-consistency implies addition of a VQG module  $G$  to generate rephrasings  $Q'$  from the image  $I$  and the predicted answer  $A'$  with an associated VQG loss  $\mathcal{L}_{vqg}(Q, Q')$ . A-consistency implies passing all the generated questions  $Q'$  to the VQA model  $F$  and an associated loss  $\mathcal{L}_{cycle}(A, A')$ . Gating implies the use of gating mechanism to filter undesirable generated questions in  $Q'$  and passing the remaining to VQA model  $F$ . Models trained with our cycle-consistent (last row in each block) framework consistently outperform baselines.

achieves an accuracy of 68.43 % on VQA v2.0 test-dev.

**Bilinear Attention Networks (BAN)** [21] <sup>4</sup> combines the idea of bilinear models and co-attention [29] between image regions and words in questions in a residual setting. Similar to [4], it uses Faster-RCNN [36] pretrained on Visual Genome [23] to extract image features. In all our experiments, for a fair comparison, we use BAN models which do not use additional training data from Visual Genome. BAN achieves the current state-of-the-art single-model accuracy of 69.64 % on VQA v2.0 test-dev without using additional training data from Visual Genome.

**Implementation Details** For all models trained with our cycle-consistent framework, we use the values  $T_{sim}=0.9$ ,  $\lambda_G=1.0$ ,  $\lambda_C=0.5$  and  $A_{iter}=5500$ . When reporting results on the validation split and VQA-Rephrasings we train on the training split and when reporting results on the test split we train on both training and validation splits of VQA v2.0. Note that we *never* explicitly train on the collected VQA-Rephrasings dataset and use it purely for evaluation purposes. We use publicly available implementations of each backbone VQA model. The hidden size of the LSTM used in VQG module is 1024 and the linear encoders used to en-

<sup>4</sup><https://github.com/jnhwkim/ban-vqa>

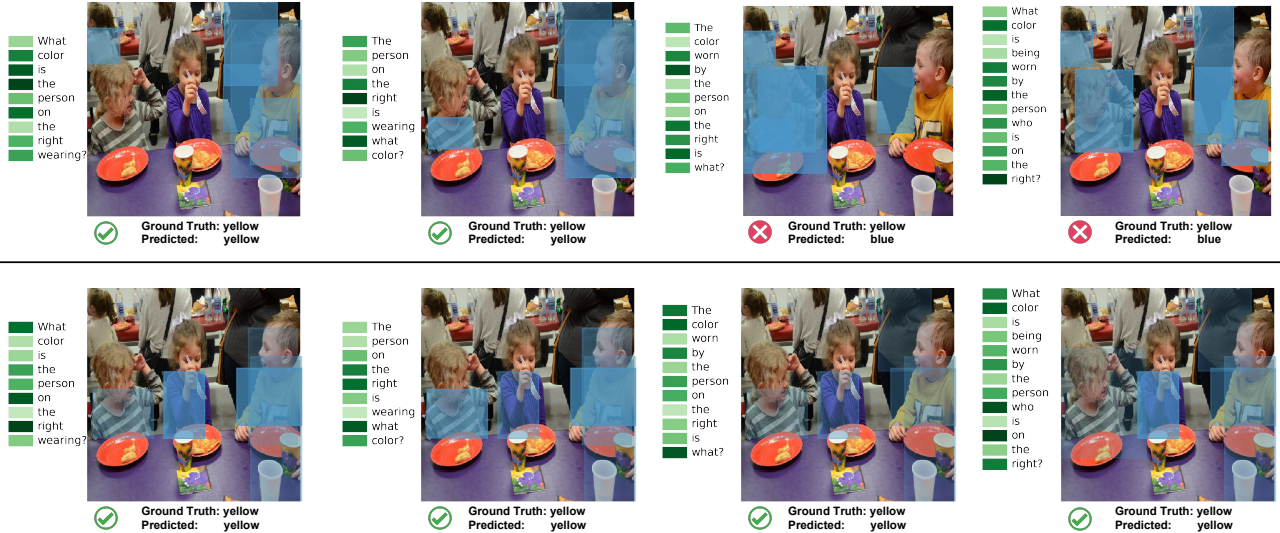


Figure 4. **Visualization of textual and image region attention across question variants:** The top row shows attention and predictions from a Pythia [46] model, the bottom row shows attention and predictions from the same Pythia model, but trained using our cycle-consistent approach. Our model attends to relevant image regions for all rephrasings and answers them correctly. The baseline Pythia counterpart, however, fails to attend over relevant image regions for some rephrasings.

Model	BLEU-1	BLEU-2	BLEU-3	BLEU-4	ROUGE-L	METEOR	CIDER
iQAN* [26]	0.582	0.467	0.385	0.320	0.617	0.276	2.222
Pythia + CC*	<b>0.708</b>	<b>0.561</b>	<b>0.438</b>	<b>0.339</b>	<b>0.627</b>	<b>0.284</b>	<b>2.301</b>
iVQA [28]	0.430	0.326	0.256	0.208	0.468	0.205	1.714
Pythia + CC	<b>0.486</b>	<b>0.368</b>	<b>0.287</b>	<b>0.226</b>	<b>0.556</b>	<b>0.225</b>	<b>1.843</b>

Table 3. **Question Generation Performance on VQA v2.0 validation set**, \* signifies results on a constrained subset as done in [26]. CC represents models trained with our approach.

code the answer and image in VQG have dimensions of 300 each. Additional details about model-specific hyperparameters can be found in Appendix E.

We measure the robustness of each of these models on our proposed VQA-Rephrasings dataset using the consensus score (Eq. 2). Table 1 shows the consensus scores at different values of  $k$  for several VQA models. We see that all models suffer significantly when measured for consistency across rephrasings. For *e.g.*, the performance of Pythia (winner of 2018 VQA challenge) is reduced to a consensus score of 39.49% at  $k = 4$ . Similar trends are observed for MUTAN, BAN and BUTD. The drop increases with increasing  $k$ , the number of rephrasings used to measure consistency. Models like BUTD, BAN and Pythia which use word-level encodings of the question suffer significant drops. It is interesting to note that even MUTAN which uses skip-thought based sentence encoding [22] suffers a drop when checked for consistency across rephrasings. We observe that BAN + CC model trained with our proposed cycle-consistent training framework consistently

outperforms its counterpart BAN and all other models at all values of  $k$ .

Fig 4 qualitatively compares the textual and visual attention (over image regions) over 4 rephrasings of a question. The top row shows attention and predictions from a Pythia model, while the bottom row shows attention and predictions from the same Pythia model, but trained using our framework. Our model attends at relevant image regions for all rephrasings and answers all of them correctly. The Pythia counterpart, however, fails to attend over relevant image regions for some rephrasings and answers those rephrasings incorrectly. This qualitatively demonstrates the robustness of models trained with our framework.

## 5.2. Visual Question Answering Performance

We now evaluate our approach and various ablations on the standard task of question answering on VQA v2.0 dataset [10]. We compare the performance of several VQA models on the validation and test-dev splits of VQA v2.0. It consists of 443,757 training, 214,354 validation and

447,793 testing questions spanning over 82,783, 40,504 and 81,434 images respectively. Table 2 shows the VQA scores of different models on validation and test-dev splits. We show that BUTD, Pythia and BAN models trained with our cycle-consistent framework outperform their corresponding baselines.

We show the impact of each component of our cycle-consistent framework by performing ablation studies on our models. We study the marginal effect of components like question consistency (Q-consistency), answer consistency (A-consistency) and gating mechanism by adding them step-by-step to the base VQA model  $F$ . Q-consistency implies addition of a VQG module  $G$  to generate rephrasings  $Q'$  from the image  $I$  and the predicted answer  $A'$  with an associated VQG loss  $\mathcal{L}_{vqg}(Q, Q')$ . As shown in Table 2, we see that addition of question consistency slightly improves performance of each VQA model. Inline with observations in [26], this shows that indeed models which can generate questions from the answer have better multi-modal understanding and in turn are better at visual question answering. A-consistency implies passing all the generated questions  $Q'$  to the VQA model  $F$  and an associated loss  $\mathcal{L}_{cycle}(A, A')$ . As seen in Table 2, we see that naively passing all the generated questions to the VQA model  $F$  leads to significant reduction in performance than the base model  $F$ . This goes in line with our earlier discussion that not all questions generated are *valid* rephrasings of the original question and hence enforcing consistency between the answers of two invalid pairs of questions naturally leads to degradation in performance. Finally we show the effect of using our gating mechanism to filter undesirable generated questions in  $Q'$  and passing the remaining to VQA model  $F$ . We see that all VQA models perform consistently better when using a gating than just using Q-consistency.

We also experimented with Pythia model configurations where the VQG model uses unattended image features (unlike the default setting which uses image features with attention from the VQA model). We found that with this configuration, our approach still shows improved performance over the baseline. However, the question generation quality is relatively poor, and the overall gain is smaller (3.58% in consistency  $CS(k=4)$  and 0.2% in VQA accuracy) compared to when using attention (8.08% and 0.5% respectively) – likely because attention helps in generating more-focused rephrasings

### 5.3. Visual Question Generation Performance

Recall that our model also includes a VQG component which generates questions conditioned on an answer and image. Since the overall performance of our framework relies highly on the performance of question generation module, we evaluate our VQG component performance as well on commonly used image captioning metrics. We compare

Model	Precision	Recall	F1
BUTD [4]	0.71	0.78	0.74
+ FP	<b>0.74</b>	<b>0.85</b>	<b>0.79</b>
BUTD + CC	0.73	0.79	0.76
+ FP	<b>0.78</b>	<b>0.83</b>	<b>0.80</b>
Pythia [46]	0.74	0.79	0.76
+ FP	<b>0.76</b>	<b>0.88</b>	<b>0.82</b>
Pythia + CC	0.77	0.81	0.77
+ FP	<b>0.82</b>	<b>0.84</b>	<b>0.83</b>

Table 4. **Failure prediction performance on VQA v2.0 validation dataset.** Each row in blocks represents a component added to the previous row. CC represents models trained with our cycle-consistent framework and FP represents models with an additional binary classification Failure Prediction submodule to predict if the predicted answer  $A'$  is correct given a question and image pair  $(Q, I)$ . For models trained without the FP module, scores are obtained by thresholding the answer confidences.

our VQG component to several answer-conditional VQG models on the VQA v2.0 dataset. We use standard image captioning metrics CIDEr [42], BLEU [34], METEOR [7] and ROUGE-L [27] as used in [28]. We compare our approach to two recently proposed visual question generation approaches. **iVQA** [28] uses a variational LSTM model trained with reinforcement learning to generate answer-specific questions for an image. Syntactic correctness, diversity and intent of the generated question are used to allocate rewards. **iQAN** [26] generates answer-specific questions by modelling question generation as a dual task of question answering and sharing parameters between question answering and question generation modules. Since iQAN can only generate a specific type of questions, for a fair comparison, we compare to iQAN only on a subset of the dataset containing questions from these specific types. As shown in Table 3, we observe that our question generation module trained with cycle-consistency consistently outperforms iVQA [28] and iQAN [26] on all metrics. A few qualitative examples of answer conditioned questions generated by our VQG model can be seen in Fig. 3(b). Additional examples can also be found in the Appendix D.

### 5.4. Failure Prediction Performance

In previous results, we show that by training models to generate and answer questions while being consistent across both tasks leads to improvement in performance and robustness. Another way of testing robustness of these models is to see if models can predict their own failures. A robust model is less confident about an incorrect answer and vice versa. Motivated by this, we seek to verify if models trained with our cycle-consistent framework can identify their own failures *i.e.* correctly identify if they’re wrong



about a prediction. To this end, we use two failure predictions schemes. First, we naively threshold the confidence of the predicted answer. All answers above a particular threshold are marked as correctly answered and vice versa. Second, we design a failure prediction binary classification module (FP), which predicts for a given image  $I$ , question  $Q$  and answer  $A'$  (predicted by the base VQA model  $F$ ), whether the predicted answer is correct for the given  $(I, Q)$  pair. The FP module uses image and answer encoders similar to those used in the question generation module (Section 3.1) and makes use of the question representation from the base VQA model as the question encoding. These encodings are concatenated and passed to a linear layer for binary classification. The FP module is trained keeping the parameters of the base VQA model frozen. In Table 4, we show the failure prediction performance of the baseline VQA models and models trained with our proposed framework. It shows that the cycle consistency framework, even *without* an explicit failure predictor module, makes the models more calibrated – more capable of detecting their own failures. In both settings: (a) when using naive confidence thresholding (not marked as “+ FP” in the Table) and (b) using a specifically designed submodule to detect failures (marked as “+ FP”), models trained with our cycle-consistent training framework are better than their corresponding baselines. We see similar improvements in detecting failures for both BUTD and Pythia models, which shows that our cycle-consistency framework is model agnostic. This also shows that not only does cycle-consistent training make models robust to linguistic variations, but also allows them to be aware of their failures.

## 6. Conclusion

In this paper, we propose a novel model-agnostic training strategy to incorporate cycle consistency in VQA models to make them robust to linguistic variations and self-aware of their failures. We also collect a large-scale dataset, VQA-Rephrasings and propose a consensus metric to measure robustness of VQA models to linguistic variations of a question. We show that models trained with our training strategy are robust to linguistic variations, and achieve state-of-the-art performance in VQA and VQG on VQA v2.0 dataset.

## References

- [1] A. Agrawal, D. Batra, and D. Parikh. Analyzing the behavior of visual question answering models. *arXiv preprint arXiv:1606.07356*, 2016.
- [2] A. Agrawal, D. Batra, D. Parikh, and A. Kembhavi. Dont just assume; look and answer: Overcoming priors for visual question answering. In *2018 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [3] H. Ali, Y. Chali, and S. A. Hasan. Automation of question generation from sentences. In *Proceedings of QG2010: The Third Workshop on Question Generation*, pages 58–67, 2010.
- [4] P. Anderson, X. He, C. Buehler, D. Teney, M. Johnson, S. Gould, and L. Zhang. 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*, 2018.
- [5] J. Andreas, M. Rohrbach, T. Darrell, and D. Klein. Learning to compose neural networks for question answering. *arXiv preprint arXiv:1601.01705*, 2016.
- [6] H. Ben-Younes, R. Cadene, M. Cord, and N. Thome. Mutan: Multimodal tucker fusion for visual question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [7] M. Denkowski and A. Lavie. Meteor universal: Language specific translation evaluation for any target language. In *Proceedings of the EACL 2014 Workshop on Statistical Machine Translation*, 2014.
- [8] A. Ettinger, S. Rao, H. Daumé III, and E. M. Bender. Towards linguistically generalizable nlp systems: A workshop and shared task. *arXiv preprint arXiv:1711.01505*, 2017.
- [9] R. Girshick, I. Radosavovic, G. Gkioxari, P. Dollár, and K. He. Detectron. <https://github.com/facebookresearch/detectron>, 2018.
- [10] 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 *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6325–6334. IEEE, 2017.
- [11] D. He, Y. Xia, T. Qin, L. Wang, N. Yu, T. Liu, and W.-Y. Ma. Dual learning for machine translation. In *Advances in Neural Information Processing Systems*, pages 820–828, 2016.
- [12] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [13] J. R. Hobbs, D. E. Appelt, J. Bear, and M. Tyson. Robust processing of real-world natural-language texts. In *Proceedings of the third conference on Applied natural language processing*, pages 186–192. Association for Computational Linguistics, 1992.
- [14] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [15] R. Hu, J. Andreas, M. Rohrbach, T. Darrell, and K. Saenko. Learning to reason: End-to-end module networks for visual question answering. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 804–813. IEEE, 2017.

- [16] M. Iyyer, J. Wieting, K. Gimpel, and L. Zettlemoyer. Adversarial example generation with syntactically controlled paraphrase networks. *arXiv preprint arXiv:1804.06059*, 2018.
- [17] U. Jain, Z. Zhang, and A. Schwing. Creativity: Generating diverse questions using variational autoencoders. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5415–5424. IEEE, 2017.
- [18] J. Johnson, B. Hariharan, L. van der Maaten, J. Hoffman, L. Fei-Fei, C. Lawrence Zitnick, and R. Girshick. Inferring and executing programs for visual reasoning. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2989–2998, 2017.
- [19] K. Kafle, M. Youssefhusien, and C. Kanan. Data augmentation for visual question answering. In *Proceedings of the 10th International Conference on Natural Language Generation*, pages 198–202, 2017.
- [20] S. Kalady, A. Elikkottil, and R. Das. Natural language question generation using syntax and keywords. In *Proceedings of QG2010: The Third Workshop on Question Generation*, 2010.
- [21] J.-H. Kim, J. Jun, and B.-T. Zhang. Bilinear Attention Networks. *arXiv preprint arXiv:1805.07932*, 2018.
- [22] R. Kiros, Y. Zhu, R. R. Salakhutdinov, R. Zemel, R. Urtasun, A. Torralba, and S. Fidler. Skip-thought vectors. In *Advances in neural information processing systems*, pages 3294–3302, 2015.
- [23] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma, et al. Visual genome: Connecting language and vision using crowd-sourced dense image annotations. *International Journal of Computer Vision*, 123(1):32–73, 2017.
- [24] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [25] Y. Li, T. Cohn, and T. Baldwin. Robust training under linguistic adversity. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, volume 2, pages 21–27, 2017.
- [26] Y. Li, N. Duan, B. Zhou, X. Chu, W. Ouyang, X. Wang, and M. Zhou. Visual question generation as dual task of visual question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
- [27] C.-Y. Lin. Rouge: A package for automatic evaluation of summaries. *Text Summarization Branches Out*, 2004.
- [28] F. Liu, T. Xiang, T. M. Hospedales, W. Yang, and C. Sun. ivqa: Inverse visual question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8611–8619, 2018.
- [29] 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*, pages 289–297, 2016.
- [30] B. McCann, N. S. Keskar, C. Xiong, and R. Socher. The natural language decathlon: Multitask learning as question answering, 2018.
- [31] I. Misra, R. Girshick, R. Fergus, M. Hebert, A. Gupta, and L. van der Maaten. Learning by asking questions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 11–20, 2018.
- [32] I. M. Mora and S. P. de la Puente. Towards automatic generation of question answer pairs from images.
- [33] N. Mostafazadeh, I. Misra, J. Devlin, M. Mitchell, X. He, and L. Vanderwende. Generating natural questions about an image. *arXiv preprint arXiv:1603.06059*, 2016.
- [34] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting on association for computational linguistics*, pages 311–318. Association for Computational Linguistics, 2002.
- [35] E. Perez, F. Strub, H. De Vries, V. Dumoulin, and A. Courville. Film: Visual reasoning with a general conditioning layer. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [36] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015.
- [37] I. V. Serban, A. García-Durán, C. Gulcehre, S. Ahn, S. Chandar, A. Courville, and Y. Bengio. Generating factoid questions with recurrent neural networks: The 30m factoid question-answer corpus. *arXiv preprint arXiv:1603.06807*, 2016.
- [38] M. Spranger, J. Suchan, and M. Bhatt. Robust natural language processing-combining reasoning, cognitive semantics and construction grammar for spatial language. *arXiv preprint arXiv:1607.05968*, 2016.
- [39] M. Stede. The search for robustness in natural language understanding. *Artificial Intelligence Review*, 6(4):383–414, 1992.
- [40] N. Sundaram, T. Brox, and K. Keutzer. Dense point trajectories by gpu-accelerated large displacement optical flow. In *European conference on computer vision*, pages 438–451. Springer, 2010.
- [41] D. Tang, N. Duan, Z. Yan, Z. Zhang, Y. Sun, S. Liu, Y. Lv, and M. Zhou. Learning to collaborate for question answering and asking. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, volume 1, pages 1564–1574, 2018.
- [42] R. Vedantam, C. Lawrence Zitnick, and D. Parikh. Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4566–4575, 2015.
- [43] Z. Wang, A. S. Lan, W. Nie, A. E. Waters, P. J. Grimaldi, and R. G. Baraniuk. Qg-net: A data-driven question generation model for educational content. In *Proceedings of the Fifth Annual ACM Conference on Learning at Scale*, pages 7:1–7:10, 2018.
- [44] X. Xu, X. Chen, C. Liu, A. Rohrbach, T. Darrell, and D. Song. Fooling vision and language models despite localization and attention mechanism. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4951–4961, 2018.

- [45] Z. Yang, X. He, J. Gao, L. Deng, and A. Smola. Stacked attention networks for image question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 21–29, 2016.
- [46] Yu Jiang\*, Vivek Natarajan\*, Xinlei Chen\*, M. Rohrbach, D. Batra, and D. Parikh. Pythia v0.1: the winning entry to the vqa challenge 2018. *arXiv preprint arXiv:1807.09956*, 2018.
- [47] P. Zhang, Y. Goyal, D. Summers-Stay, D. Batra, and D. Parikh. Yin and yang: Balancing and answering binary visual questions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5014–5022, 2016.
- [48] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Computer Vision (ICCV), 2017 IEEE International Conference on*, 2017.

The appendix is organized as follows:

- Section A covers information about the dataset collection pipeline, user interface and provides some dataset statistics.
- Section B shows qualitative examples of how attention over image regions varies for VQA models when different rephrasings of the same question are used as input.
- Section C describes an attention based consistency strategy that we experimented with, but did not improve performance (and so was not a part of our final model presented in the paper).
- Section D shows qualitative examples of answer conditioned questions generated by our VQG module.
- Section E lists the hyperparameters used for each base VQA model.

## A. Dataset Details

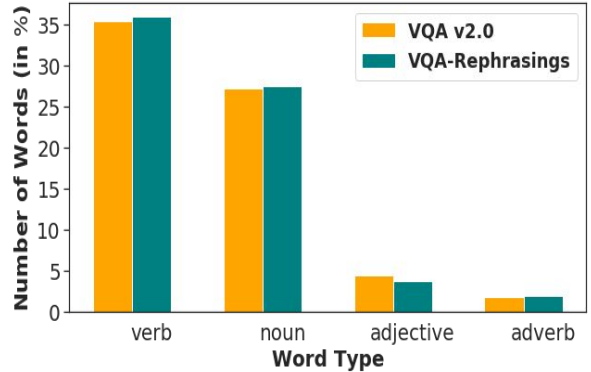
**Statistics.** Fig 5(a) shows the number of words (in percentage) belonging to different Parts-of-Speech tags. The distributions follow almost similar trends in VQA-Rephrasings and VQA v2.0. This shows that the rephrasings are not obtained by merely adding more adjectives or adverbs in the original question. Fig 5(b) shows the number of questions (in percentage) with varying lengths. The average length of questions in VQA-Rephrasings is 7.15 which is slightly higher than the average length in VQA v2.0, which is 6.32.

**Interface.** We used a simplistic web interface to collect rephrasings from human annotators. The interface provided three examples of invalid rephrasings and their corresponding explanations to help human annotators understand the task better. We A/B tested with 50 questions using all 4 combinations of:

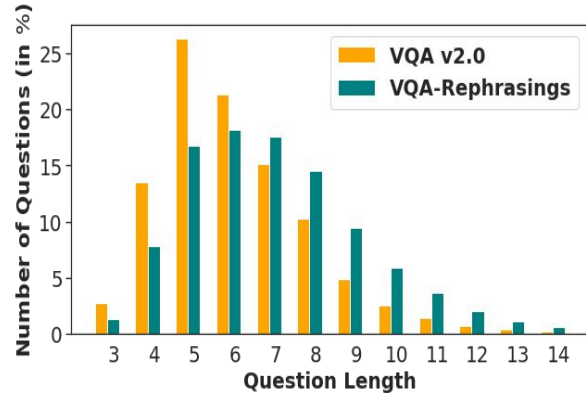
- Showing both valid and invalid rephrasing examples and explanations.
- Showing only valid and no invalid rephrasing examples and explanations.
- Showing none of valid and invalid rephrasing examples and explanations.
- Showing no valid and only invalid rephrasing examples and explanations.

We found (via manual inspection) that the last setup provided higher quality data, and used that as our final interface choice.

**Examples.** Fig 6 shows several qualitative examples from the VQA-Rephrasings dataset. We see that the rephrasings maintain the intent of the original question while varying linguistically.



(a)



(b)

Figure 5. **Dataset Statistics.** (a) Shows the number of words (in percentage) belonging to different Parts-of-Speech tags. The distributions follow similar trends in VQA-Rephrasings and VQA v2.0. (b) Shows the number of questions (in percentage) with varying lengths. The average length of questions in VQA-Rephrasings is 7.15 which is slightly higher than the average length in VQA v2.0, which is 6.32.

## B. Attention Analysis

Fig 7 qualitatively compares the textual and visual attention (over image regions) for rephrasings of a question. Each row compares predicted answers and attention from a baseline Pythia [46] model and the same Pythia model trained with our framework (Pythia + CC), using two question rephrasings. First and third row shows the outputs of a Pythia model (baseline) and second and forth row shows the output of a Pythia model (baseline + CC) trained with our framework. We see that in most examples, the attention over image regions doesn't vary across rephrasings for models trained with our framework (and the model answers the questions correctly). However for the baseline model, one can see that minor linguistic changes in the question can result in completely different answers (Row 2, Columns

1 and 3). This qualitatively demonstrates the robustness of models trained with our framework. Since the baseline Pythia model doesn't include a counting module, it doesn't perform well on questions requiring counting. As a result we see that both the baseline and its cycle-consistent counterpart perform poorly on counting questions (Row 5, Columns 1 through 4).

### C. Attention Consistency

Intuitively, it seems like training the VQA model to attend over the same image regions for different rephrasings of a question should improve the robustness of the model. We tried to enforce this in our cycle-consistent framework using an additional attention consistency loss.

Recall that for a given image  $I$ , question  $Q$  and answer  $A$ , our model consists of a VQA model  $F$  which takes  $(Q, I)$  as an input and uses the question to attend over image regions with attention  $\gamma_Q$  and predicts an answer  $A'$ . We also have a VQG model  $G$  which uses the predicted answer  $A'$  and image  $I$  to generate a question  $Q'$ . Intuitively, the VQA model should attend over the same image regions when answering  $Q'$ . In other words, the attention over image regions  $\gamma_{Q'}$  used by the VQA model to answer  $Q'$  should be close to the  $\gamma_Q$ . We added an additional attention consistency loss to the total loss which reduces the  $L_2$  norm between these two attentions.

However, we found that this leads to reduction in model performance. Specifically, this reduces the performance of a cycle consistent Pythia model by 1.34% VQA accuracy when evaluated on the VQA v2.0 validation split (training on train split only).

We suspect one reason why enforcing attention consistency across rephrasings reduces performance is perhaps because minimizing a large number of diverse losses (cross entropy losses  $\mathcal{L}_F$  and  $\mathcal{L}_{cycle}$  for VQA, sequence generation loss  $\mathcal{L}_G$  for VQG and mean squared loss  $\mathcal{L}_{attention}$  for attention consistency) is a hard problem to optimize. Concretely identifying why enforcing attention consistency across question rephrasings hurts performance is currently under investigation and is part of future work. We find naively matching attentions across question rephrasings is not effective in current settings and therefore do not include this in the final model.

### D. Question Generation

Fig 8 shows qualitative examples of answer conditioned questions generated by our VQG model. Our VQG model is able to correctly generate answer conditioned questions for a wide range of answers ranging from numbers, to colors and even yes/no.

### E. Hyperparameters

We use the default hyperparameters as described in publicly available implementations of MUTAN [6], BUTD [4], Pythia [46] and BAN [21]. When using these models as base VQA models to train cycle consistent variants of them, we use the same parameters for the VQA model. For the the VQG model we use  $T_{sim}=0.9$ ,  $\lambda_G=1.0$ ,  $\lambda_C=0.5$  and  $A_{iter}=5500$ . While some models use adaptive learning rates for their base VQA models, the VQG model is always trained with a fixed learning rate of 0.0005. In case of BAN and Pythia, we also clip the gradients whose  $L_2$  norm is greater than 0.25.



	<ul style="list-style-type: none"> <li>Where is the nike sign?</li> <li>Where can I find the nike sign?</li> <li>Where is the nike sign located?</li> <li>What is the location of the nike sign?</li> </ul>		<ul style="list-style-type: none"> <li>What do the orange words say?</li> <li>What does the orange words read?</li> <li>What is written in orange text?</li> <li>What does it say in orange?</li> </ul>
	<ul style="list-style-type: none"> <li>Is the horse running?</li> <li>Does the horse appear to be running?</li> <li>Does it look like the horse is running?</li> <li>Is the horse in a running motion?</li> </ul>		<ul style="list-style-type: none"> <li>What kind of food is the green items?</li> <li>What is the green food?</li> <li>What do you call the green food pictured?</li> <li>What is the green food item known as?</li> </ul>
	<ul style="list-style-type: none"> <li>Is this in a cold climate?</li> <li>Is the climate here cold?</li> <li>Is a cold climate shown here?</li> <li>Is the climate here frigid?</li> </ul>		<ul style="list-style-type: none"> <li>Are all the ducks swimming?</li> <li>Is every duck swimming?</li> <li>Is each duck swimming?</li> <li>Are all of the ducks on the water swimming?</li> </ul>
	<ul style="list-style-type: none"> <li>How high is the plane in the sky?</li> <li>What altitude is the plane flying at?</li> <li>How high up in the air is the plane?</li> <li>Do you know the plane's current altitude?</li> </ul>		<ul style="list-style-type: none"> <li>Are the planes planning to land soon?</li> <li>Do you think the planes will land shortly?</li> <li>Are the planes going to land soon?</li> <li>Do you anticipate the planes to land soon?</li> </ul>
	<ul style="list-style-type: none"> <li>Are the children related?</li> <li>Are the kids related to each other?</li> <li>Are the children relatives of each other?</li> <li>Do those children come from the same family?</li> </ul>		<ul style="list-style-type: none"> <li>What game are they playing?</li> <li>What game is everyone participating in?</li> <li>What is everyone playing?</li> <li>What is the game called that the people are playing?</li> </ul>
	<ul style="list-style-type: none"> <li>Would a vegetarian eat this meal?</li> <li>If you were a vegetarian would you eat this meal?</li> <li>Is this a meal a vegetarian would eat?</li> <li>Would this be a meal a vegetarian would eat?</li> </ul>		<ul style="list-style-type: none"> <li>Was this pizza homemade?</li> <li>Is this a homemade pizza?</li> <li>Was the pizza made at home?</li> <li>Is the pizza considered to be homemade?</li> </ul>
	<ul style="list-style-type: none"> <li>Was this food cooked in a oven?</li> <li>Is the oven what the food was cooked in?</li> <li>Was the food prepared in an oven?</li> <li>Was the oven used to cook the food?</li> </ul>		<ul style="list-style-type: none"> <li>Is this animal dangerous?</li> <li>Should you be afraid of this animal?</li> <li>Is this a dangerous animal?</li> <li>Is this animal threatening?</li> </ul>
	<ul style="list-style-type: none"> <li>Is there a white horse running?</li> <li>Is a white horse running in the picture?</li> <li>Is there a horse that is white colored running?</li> <li>Can you see a white colored horse running?</li> </ul>		<ul style="list-style-type: none"> <li>Is this vegetable better cooked?</li> <li>Is cooked better than raw for this vegetable?</li> <li>Is this vegetable preferred cooked?</li> <li>Would you say this vegetable is better if cooked?</li> </ul>
	<ul style="list-style-type: none"> <li>How many more hours until midnight?</li> <li>Midnight is in how many hours?</li> <li>What's the number of hours until midnight?</li> <li>How many hours to go until it's midnight?</li> </ul>		<ul style="list-style-type: none"> <li>What is the occasion that this photo depicts?</li> <li>What occasion is this photo showing?</li> <li>What is the occasion that is shown in this photo?</li> <li>Which occasion does this picture depict?</li> </ul>
	<ul style="list-style-type: none"> <li>Which sign is for a fast food company?</li> <li>What fast food company is this sign for?</li> <li>The sign featured is for what fast food company?</li> <li>What fast food company has this sign?</li> </ul>		<ul style="list-style-type: none"> <li>Are there any spices on the pizza?</li> <li>Does the pizza have spices on it?</li> <li>Is the pizza garnished with any spices?</li> <li>Are there some sort of spices on the pizza?</li> </ul>
	<ul style="list-style-type: none"> <li>Is this a low calorie meal?</li> <li>Is this meal healthy?</li> <li>Is this a healthy meal?</li> <li>Does the food look like a low calorie meal?</li> </ul>		<ul style="list-style-type: none"> <li>What is this person doing?</li> <li>What activity is this person participating in?</li> <li>How would you describe the person's activity?</li> <li>What activity is the person engaged in?</li> </ul>

Figure 6. Examples from our VQA-Rephrasings dataset. The first question (shown in gray) in each block is the original question from VQA v2.0 validation set, the questions that follow (shown in black) are rephrasings collected in VQA-Rephrasings.

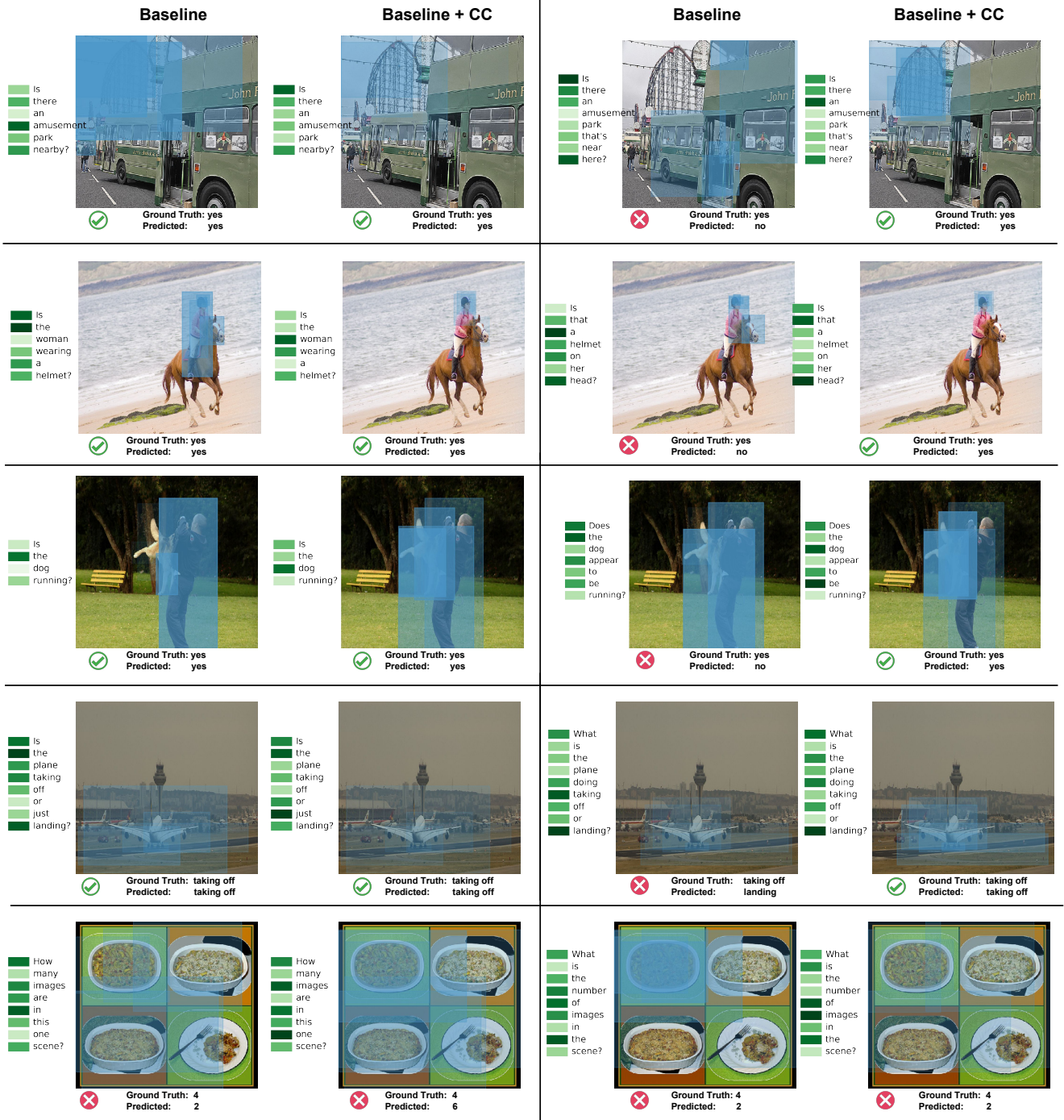


Figure 7. Visualization of textual and image region attention for different question variants: Each row compares answers predicted and attention for two question rephrasings using a baseline Pythia [46] model and the same Pythia model trained with our framework (Pythia + CC). Higher opaqueness in highlighted regions represents higher attention. First and third rows show the output of a Pythia model (baseline) and second and forth rows show the output of a Pythia model (baseline + CC) trained with our framework. As one can see, in most examples, the attention over image regions doesn't vary much for models trained with our framework. However for the baseline model, one can see that by very minor linguistic changes in the question it is possible to predict completely different answers (Row 2, Columns 1 and 3). These examples qualitatively demonstrate the robustness of models trained with our framework.





- **One**
  - How many chairs are in the room?
  - How many beds are present?
- **No**
  - Is this a hotel room?
  - Is there a person in the room?
- **Yes**
  - Is the bed made?
  - Is the bed white in color?
  - Is the wall white?
  - Does the desk look messy?



- **Yes**
  - Is this a winter a scene?
  - Are there trees in the background?
  - Is the man in air?
  - Is the man snowboarding?
- **No**
  - Is the man wearing goggles?
  - Is there snow on the trees?
- **Red**
  - What color jacket is the man wearing?



- **No**
  - Is the bus moving?
  - Is the man in picture on the right side of the bus?
- **Concrete**
  - What is the sidewalk made of?
- **Gray**
  - What is the color of the bus?
  - What color is the vehicle in the picture?
- **Daytime**
  - Is it daytime or nighttime?



- **Blue**
  - What is the color of the suitcase?
- **Red**
  - What color is the door?
  - What is the color of the door?
- **No**
  - Is the man carrying a backpack?
- **One**
  - How many suitcases are there?
  - How many suitcases can be seen?



- **Yes**
  - Are there any chairs in the picture?
  - Are there flowers in the picture?
  - Are the flowers in a garden
- **Seven**
  - How many chairs are there in the picture?
- **Grass**
  - What is the object on the left?



- **Apple**
  - What fruit is shown in the picture?
  - What kind of fruit is that?
  - Which fruit is shown in the picture?
- **Red**
  - What color is the fruit shown?
  - What color are the apples?
- **10**
  - How many birds can be seen in the picture?

Figure 8. Qualitative examples of answer conditioned question generation by our VQG module.