

The Promise of Premise: Harnessing Question Premises in Visual Question Answering

Aroma Mahendru^{*,1} Viraj Prabhu^{*,1} Akrit Mohapatra^{*,1} Dhruv Batra² Stefan Lee¹
¹Virginia Tech ²Georgia Institute of Technology
{maroma, virajp, akrit}@vt.edu, dbatra@gatech.edu, steflee@vt.edu

Abstract

In this paper, we make a simple observation that questions about images often contain *premises* – objects and relationships implied by the question – and that reasoning about premises can help Visual Question Answering (VQA) models respond more intelligently to irrelevant or previously unseen questions.

When presented with a question that is irrelevant to an image, state-of-the-art VQA models will still answer purely based on learned language biases, resulting in nonsensical or even misleading answers. We note that a visual question is irrelevant to an image if at least one of its premises is false (*i.e.* not depicted in the image). We leverage this observation to construct a dataset for Question Relevance Prediction and Explanation (QRPE) by searching for false premises. We train novel question relevance detection models and show that models that reason about premises consistently outperform models that do not.

We also find that forcing standard VQA models to reason about premises during training can lead to improvements on tasks requiring compositional reasoning.

1 Introduction

The task of providing natural language answers to free-form questions about an image – *i.e.* Visual Question Answering (VQA) – has received substantial attention in the past few years (Malinowski and Fritz, 2014; Antol et al., 2015; Malinowski et al., 2015; Zitnick et al., 2016; Kim et al., 2016; Wu et al., 2016; Lu et al., 2016; Andreas et al.,

What brand of **racket** is the **man** holding ?

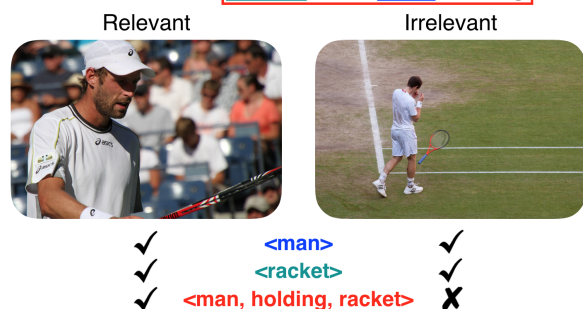


Figure 1: Questions asked about images often contain ‘premises’ that imply visual semantics. From the above question, we can infer that a relevant image must contain a man, a racket, and that the man must be holding the racket. We extract these premises from visually grounded questions and use them to construct a new dataset and models for question relevance prediction. We also find that augmenting standard VQA training with simple premise-based questions results in improvements on tasks requiring compositional reasoning.

2016; Lu et al., 2017) and has quickly become a popular problem area. Despite significant progress on VQA benchmarks (Antol et al., 2015), current models still present a number of unintelligent and problematic tendencies.

When faced with questions that are irrelevant or not applicable for an image, current ‘forced choice’ models will still produce an answer. For example, given an image of a dog and a query “What color is the bird?”, standard VQA models might answer “Red” confidently, based solely on language biases in the training set (*i.e.* an overabundance of the word “red”). In these cases, the predicted answers are senseless at best and misleading at worst, with either case posing serious problems for real-world applications. Like Ray et al. (2016), we argue that practical VQA systems must be able to identify and explain irrelevant questions. For instance, a more intelligent VQA model with this capability might answer “There is no bird in the image” for this example.

*Denotes equal contribution.

Premises. In this paper, we show that question *premises* - *i.e.* objects and relationships implied by a question - can enable VQA models to respond more intelligently to irrelevant or previously unseen questions. We develop a premise extraction pipeline based on SPICE (Anderson et al., 2016) and demonstrate how these premises can be used to improve modern VQA models in the face of irrelevant or previously unseen questions.

Concretely, we define premises as facts implied by the language of questions, for example the question “*What brand of racket is the man holding?*” shown in Fig. 1 implies the existence of a man, a racket, and that the man is holding the racket. For visually grounded questions (*i.e.* those asked about a particular image) these premises imply visual qualities, including the presence of objects as well as their attributes and relationships.

Broadly speaking, we explore the usefulness of premises in two settings – when visual questions are known to be relevant to the images they are asked on (*e.g.* in the VQA dataset) and in real-life situations where such an assumption cannot be made (*e.g.* when generated by visually impaired users). In the former case, we show that knowing that a question is relevant allows us to perform data augmentation by creating additional simple question-answer pairs using the premises of source questions. In the latter case, we show that explicitly reasoning about premises provides an effective and interpretable way of determining whether a question is relevant to an image.

Irrelevant Question Detection. We consider a question to be relevant to an image if all of the question’s premises apply to the corresponding image, that is to say all objects, attributes, and interactions implied by the question are depicted in the image. We refer to premises that apply for a given image as true premises and those that do not apply as false premises. In order to train and evaluate models for this task, we curate a new irrelevant question detection dataset which we call the Question Relevance Prediction and Explanation (QRPE) dataset. QRPE is automatically curated from annotations already present in existing datasets, requiring no additional labeling.

We collect the QRPE dataset by taking each image-question pair in the VQA dataset (Antol et al., 2015) and finding the most visually similar other image for which exactly one of the question premises is false. In this way, we collect tu-

ples consisting of two images, a question, and a premise where the question is relevant for one image and not for the other due to the premise being false.

For context, the only other existing irrelevant question detection dataset (Ray et al., 2016) collected irrelevant question-image pairs by human verification of random pairs. In comparison, QRPE is substantially larger, balanced between irrelevant and relevant examples, and presents a considerably more difficult task due to the closeness of the image pairs both visually and with respect to question premises.

We train novel models for irrelevant question detection on the QRPE dataset and compare to existing methods. In these experiments, we show that models that explicitly reason about question premises consistently outperform baseline models that do not.

VQA Data Augmentation. Finally, we also introduce an approach to generate simple, templated question-answer pairs about elementary concepts from premises of complex training questions. In initial experiments, we show that adding these simple question-answer pairs to VQA training data can improve performance on tasks requiring compositional reasoning. These simple questions improve training by bringing implicit training concepts “to the surface”, *i.e.* introducing direct supervision of important implicit concepts by transforming them to simple training pairs.

2 Related Work

Visual Question Answering: Starting from simple bag-of-words and CNN+LSTM models (Antol et al., 2015), VQA architectures have seen considerable innovation. Many top-performing models integrate attention mechanisms (over the image, the question, or both) to focus on important structures (Fukui et al., 2016; Lu et al., 2016, 2017), and some have been designed with compositionality in mind (Andreas et al., 2016; Hendricks et al., 2016). However, improving compositionality or performance through data augmentation remains a largely unstudied area.

Some other recent work has developed models which produce natural language explanations for their outputs (Park et al., 2016; Wang et al., 2016), but there has not been work on generating explanations for irrelevant questions or false premises.

Question Relevance: Most related to our work is that of Ray et al. (2016), which introduced the task of irrelevant question detection for VQA. To evaluate on this task, they created the Visual True and False Question (VTFQ) dataset by pairing VQA questions with random VQA images and having human annotators verify whether or not the question was relevant. As a result, many of the irrelevant image-question pairs exhibit a complete mismatch of image and question content. Our Question Relevance Prediction and Explanation (QRPE) dataset on the other hand is collected such that irrelevant images for each question closely resemble the source image both visually and semantically. We also provide premise-level annotations which can be used to develop models that not only decide whether a question is relevant, but also provide explanations for *why* that is the case.

Semantic Tuple Extraction: Extracting structured facts in the form of semantic tuples from text is a well studied problem (Schuster et al., 2015; Anderson et al., 2016; Elhoseiny et al., 2016); however, recent work has begun extending these techniques to visual domains (Xu et al., 2017; Johnson et al., 2015). Additionally, the Visual Genome (Krishna et al., 2016) dataset contains dense image annotations for objects and their attributes and relationships. However, we are the first to consider these facts to reason about question relevancy and compositionality in VQA.

3 Extracting Premises of a Question

In Section 1, we introduced the concept of premises and how they can be used. We now formalize this concept and explain how premises can be extracted from questions.

We define question premises as facts implied about an image from a question asked about it, which we represent as tuples. Returning to our running example question “*What brand of racket is the man holding?*”, we can express these premises as the tuples ‘<man>’, ‘<racket>’, and ‘<man, holding, racket>’ respectively. We categorize these tuples into three groups based on their complexity. First-order premises representing the presence of objects (‘<man>’, ‘<cat>’, ‘<sky>’), second-order premises capturing the attributes of objects (‘<man, tall>’, ‘<car, moving>’), and third-order premises containing interactions between objects (e.g. ‘<man, kicking, ball>’, ‘<cat, above, car>’).

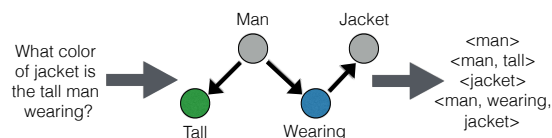


Figure 2: **Premise Extraction Pipeline.** Objects (gray), attributes (green), and relations (blue) scene graph nodes are converted into 1st, 2nd, and 3rd order premises respectively.

Premise Extraction: To extract premises from questions, we use the semantic tuple extraction pipeline used in the SPICE metric (Anderson et al., 2016). Originally defined as a metric for image captioning, SPICE transforms a sentence into a scene graph using the Stanford Scene Graph Parser (Schuster et al., 2015) and then extracts semantic tuples from this representation. Fig. 2 shows this process for a sample question. The question is represented as a graph of objects, attributes, and relationships from which first, second, and third order premises are extracted respectively. As this pipeline was originally designed for descriptive captions rather than questions, we found a number of minor modifications helpful in extracting quality question premises, including disabling pronoun resolution, verb lemmatization and METEOR-based Synset matching. We will release our premise extraction code publicly to encourage reproducibility.

While this extraction process typically produces high quality premise tuples, there are some sources of noise which must be filtered out. The SPICE process occasionally produces duplicate nodes or object nodes not linked to nouns in the question, which we filter out. We also remove premises containing words like photo, image, *etc.* that refer to the image rather than its content.

A more nuanced source of erroneous premises comes from the ambiguity in existential questions, *i.e.* those about the existence of certain image content. For example, while the question “*Is the little girl moving?*” contains the premise ‘<girl, little>’, it is unclear without the answer whether ‘<girl, moving>’ is also a premise. Similarly, for the question “*How many giraffes are in the image?*”, ‘<giraffe, many>’ cannot be considered a premise as there may be 0 giraffes in the image. To avoid introducing false premises, we filter out existential and counting questions.



Figure 3: **Some Examples from QRPE Dataset.** For a given question Q and a relevant image I^+ , we find an irrelevant image I^- for which exactly one premise P of the question is false. If there are multiple such candidates, we select the candidate most visually most similar to I^+ . As can be seen from these examples, the QRPE dataset is very challenging, with only minor visual and semantic differences separating the relevant and irrelevant images.

4 Question Relevance Prediction and Explanation (QRPE) Dataset

As discussed in Section 1, modern VQA models fail to differentiate between relevant and irrelevant questions, answering either with confidence. This behavior is detrimental to the real world application of VQA systems. In this section, we curate a new dataset for question relevance in VQA which we call the Question Relevance Prediction and Explanation (QRPE) dataset. We plan to release QRPE publicly to help future efforts.

In order to train and evaluate models for irrelevant question detection, we would like to create a dataset of tuples (I^+, Q, P, I^-) comprised of a natural language question Q , an image I^+ for which Q is relevant, and an image I^- for which Q is irrelevant because premise P is false. While it is not required to collect both a relevant and irrelevant image for each question, we argue that doing so is a simple way to balance the dataset and it ensures that biases against rarer questions (which would be irrelevant for most images) cannot be exploited to inflate performance.

We base our dataset on the existing VQA corpus (Antol et al., 2015), taking the human-generated (and therefore relevant) image-question pairs from VQA as I^+ and Q . As previously discussed, we can define the relevancy of a question in terms of the validity of its premises for an image, so we extract premises from each question Q and must find a suitable irrelevant image I^- . However, there are certainly many images for which one or more of Q 's premises are false and an important design decision is then how to select I^- from this set.

To ensure our dataset is as realistic and challenging as possible, we consider irrelevant images which only have a single false question premise under Q which we denote P . For example, the question “*Is the big red dog old?*” could be matched with an image containing a big, white dog or a small red dog, but not a small white dog. In this way, we ensure that image content is semantically appropriate for the question topic but not quite relevant. Additionally, this provides each irrelevant image with an explanation for why the question does not apply.

Furthermore, we sort this subset of irrelevant image by their visual distance to the source image I^+ based on image encodings from a VGGNet (Simonyan and Zisserman, 2014) pretrained on ImageNet (Russakovsky et al., 2012). This ensures that the relevant and irrelevant images are visually similar and act as difficult examples.

A major difficulty with our proposed data collection process is how to verify whether a premise is true or false for any given image in order to identify irrelevant images. We detail dataset construction and our approach for this problem in the following section.

4.1 Dataset Construction

We curate our QRPE dataset automatically from existing annotations in COCO (Lin et al., 2014) and Visual Genome (Krishna et al., 2016). COCO is a set of over 300,000 images annotated with object segmentations and presence information for 80 classes as well as text descriptions of image content. Visual Genome builds on this dataset, providing more detailed object, attribute, and rela-

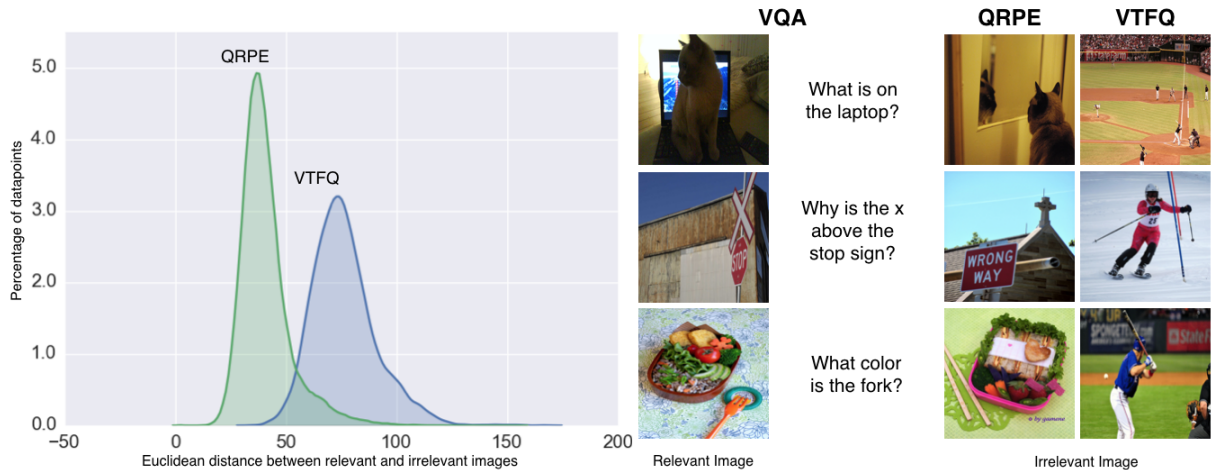


Figure 4: **A comparison of the QRPE and VTFQ Datasets.** On the left, we plot the Euclidean distance between VGGNet-*fc7* features extracted from each relevant-irrelevant image pair for each dataset. Note that VTFQ has significantly higher visual distances. On the right, we show some qualitative examples of irrelevant images for questions that occur in both datasets. VTFQ images are significantly less related to the source image and question than in our dataset.

tionship annotations for over 100,000 COCO images. We make use of these data sources to extract first and second order premises from VQA questions which are also based on COCO images.

For first order premises (*i.e.* existential premises), we consider only the 80 classes present in COCO (Lin et al., 2014). As VQA and COCO share the same images, we can easily determine if a first order premise is true or false for a candidate irrelevant image simply by checking for the absence of the appropriate class annotation.

For second order premises (*i.e.* attributed objects), we rely on Visual Genome (Krishna et al., 2016) annotations for object and attribute labels. Unlike in COCO, the lack of a particular object label in an image for Visual Genome does not necessarily indicate that the object is not present, both due to annotation noise and the use of multiple synonyms for objects by human labelers. As a consequence, we restrict the set of candidate irrelevant images to those which contain a matching object to the question premise but a different attribute. Without further restriction, the selected irrelevant attributes do not tend to be mutually exclusive with the source attribute (*i.e.* matching ‘<dog, old>’ and ‘<dog, red>’). To correct this and ensure a false premise, we further restrict the set to attributes which are antonyms (*e.g.* ‘<young>’ for source attribute ‘<old>’) or taxonomic sister terms (*e.g.* ‘<green>’ for source attribute ‘<red>’) of the original premise attribute. We also experimented with third order premises;

however, the lack of a corresponding sense of mutual exclusion for verbs and the sparsity of <object, relationship, object> premises made finding non-trivial irrelevant images difficult.

To recap, our data collection approach is to take each image-question pair in the VQA dataset and extract its first and second order question premises. For each premise, we find all images which lack only this premise and rank them by their visual distance. The closest of these is kept as the irrelevant image for each image-question pair.

4.2 Exploring the Dataset

Fig. 3 shows sample (I^+, Q, P, I^-) tuples from our dataset. These examples illustrate the difficulty of our dataset. For instance, the images in the second column differ only by the presence of the water bottle and images in the fourth column are differentiated by the color of the devices. Both of these are fine details of the image content.

The QRPE dataset contains 53,911 (I^+, Q, P, I^-) tuples generated from as many premises. In total, it contains 1530 unique premises and 28,853 unique questions. Among the 53,911 premises, 3876 are second-order, attributed object premises while the remaining 50,035 are first-order object/scene premises. We divide our dataset into two parts – a training set with 35,486 tuples that are generated from the VQA training set and a validation set with 18,425 tuples generated from the VQA validation set.

Manual Validation. We also manually validated 1000 randomly selected (I^+, Q, P, I^-) tuples from our dataset. We noted that 99.10% of the premises P were valid (*i.e.* implied by the question) in I^+ and 97.3% were false for the negative image I^- . This demonstrates the high reliability of our automated annotation pipeline.

4.3 Comparison to VTFQ

We contrast our approach to the VTFQ dataset of Ray et al. (2016). As discussed prior, VTFQ was collected by selecting a random question and image from the VQA set and asking human annotators to report if the question was relevant, producing a pair. This approach results in irrelevant image-question pairs that are unambiguously unrelated, with the visual content of the image having nothing at all to do with the question or its source image from VQA.

To quantify this effect and compare to QRPE, we pair each irrelevant image-question pair (I^-, Q) from VTFQ with a relevant image from the VQA dataset. Specifically, we find the nearest neighbor question Q^{nn} in the VQA dataset to Q based on an average of the word2vec (Mikolov et al., 2013) embedding of each word, and select the image on which Q^{nn} was asked as I^+ to form (I^+, Q, P, I^-) tuples like in our proposed dataset.

In Fig. 4, we present a quantitative and qualitative comparison of the two datasets based on these tuples. On the left side of the figure, we plot the distributions of Euclidean distance between the fc7 features of each (I^+, I^-) pair in both datasets. We find that the mean distance in the VTFQ dataset is nearly twice that of our QRPE dataset, indicating that irrelevant images in VTFQ are less visually related to source images though we do note the distribution of distances in both datasets is long tailed.

On the right side of Fig. 4, we also provide qualitative examples of questions that occur in both datasets. The example on the last row is perhaps most striking. The source question is asking the color of a fork and the relevant image shows an overhead view of a meal with an orange fork set nearby. The irrelevant image in QRPE is a similar image of food, but with chopsticks! Conversely, the image from VTFQ is a man playing baseball.

5 Question Relevance Detection

In this section, we introduce a simple baseline for irrelevant question detection on the QRPE dataset and demonstrate that explicitly reasoning about premises improves performance for both our new model and existing methods. More formally, we consider the binary classification task of predicting if a question Q_i from an image-question pair (I_i, Q_i) is relevant to image I_i .




A Simple Premise-Aware Model. Like the standard VQA task, question relevance detection also requires making a prediction based on an encoded image and question. With this in mind, we begin with a straight-forward approach based on the Deeper LSTM VQA model architecture of Antol et al. (2015). This model encodes the image I via a VGGNet and the question Q with an LSTM over one-hot word encodings. The concatenation of these embeddings are input to a multi-layer perceptron. We fine-tune this model for the binary question relevance detection task starting from a model pretrained on the VQA task. We denote this model as VQA-Bin.

We extend the VQA-Bin model to explicitly reason about premises. We extract first and second order premises from the question Q and encode them as two concatenated one-hot vectors. We add an additional LSTM to encode the premises and concatenate this added feature to the image and question feature. We refer to this premise-aware model as VQA-Bin-Premise.

Attention Models. We also extend the attention based Hierarchical Co-Attention VQA model of Lu et al. (2016) for the task of question relevance in a way similar to Deeper LSTM model. We call this model HieCoAtt-Bin. The corresponding premise-aware model is referred to as HieCoAtt-Bin-Prem.

Existing Methods. We compare our approaches with the best performing model of Ray et al. (2016). This model (which we denote QC-Sim) uses a pretrained captioning model to automatically provide natural language image descriptions and reasons about relevance based on a learned similarity between the question and image caption.

Specifically, the approach uses NeuralTalk2 (Karpathy and Li, 2015) trained on the MS COCO dataset (Lin et al., 2014) to generate a caption for each image. Both the caption and question are

Question	Premise	Valid	Explanation
 What color is the cat's tie?	<cat> <tie>	✗ ✗	There is no cat There is no tie
 What kind of building is the large white building?	<building, large> <building, white>	✓ ✓	There is a large building There is a white building
 What is the polar bear in?	<bear, polar>	✗	There is no polar bear




Question	Premise	Valid	Explanation
 What color is the car behind the bus?	<car> <bus>	✓ ✓	There is a car There is a bus
 What type of birds are on the bench?	<bird> <bench>	✗ ✗	There is no bird There is no bench
 Where is the small elephant?	<elephant> <elephant, small>	✓ ✓	There is an elephant There is a small elephant

Figure 5: **Question relevance explanation:** We provide selected examples of predictions from the False Premise Detection model (FPD) on the QRPE test set. Reasoning about premises presents the opportunity to produce natural language statements indicating *why* a question is irrelevant to an image, by pointing to the premise that is invalid.

Models	Overall	First Order	Second Order
VQA-Bin	66.50	67.36	53.00
VQA-Bin-Prem	66.77	67.04	54.38
HieCoAtt-Bin	70.74	71.35	61.54
HieCoAtt-Bin-Prem	73.34	73.97	60.35
QC-Sim	74.35	75.82	55.12
PC-Sim	75.05	76.47	56.04
QPC-Sim	75.31	76.67	55.95

Table 1: Accuracy of Question Relevance models on the QRPE test set. We find that premise-aware models consistently outperform alternative models.

embedded as a fixed length vector through an encoding LSTM (with words being represented as word2vec (Mikolov et al., 2013) vectors). These question and caption embeddings are concatenated and fed to a multilayer perceptron to predict relevance. We consider two additional versions of this approach that consider only premise-caption similarity (PC-Sim) and question-premise-caption similarities (QPC-Sim).

Results. We train each model on the QRPE train split and report results on the test set in Table 1. As the dataset is balanced in the label space, random accuracy stands at 50%. We find that the simple VQA-Bin model achieves 66.5% accuracy while the attention based model HieCoAtt-Bin attains 70.74% accuracy. Surprisingly, the caption-similarity based QC-Sim model significantly outperforms these baseline, obtaining an accuracy of 74.35% while only reasoning about relevancy from textual descriptions of images. We note that the caption similarity based approaches use a large amount of outside data during pretraining of the captioning model and the word2vec embeddings, which may have contributed to the effectiveness of these methods.

Most interestingly, we find that the addition of extracted premise representations consistently improves performance of base models. VQA-Bin-Prem, HieCoAtt-Bin-Prem, PC-Sim, and QPC-Sim outperform their no-

premise information counterparts, with QPC-Sim being the overall best performing approach at 75.31% accuracy. This is especially interesting given that the models *already* have access to the question from which the premises were extracted. This result seems to imply there is value in explicitly isolating premises from sentence grammar.

We further divide our test set into two splits consisting of (Q, I) pairs created by either falsifying first-order and second-order premises. We find that all our models perform significantly better on the first-order split. We hypothesize that the significant diversity in visual representations of attributed objects and comparatively fewer examples for each type makes it more difficult to learn subtle differences for second-order premises.

5.1 Question Relevance Explanation

In addition to identifying whether a question is irrelevant to an image, being able to indicate *why* carries significant real-world utility. From an interpretability perspective, reporting which premise is false is more informative than simply answering the question in the negative, as it can help to correct the questioner’s misconception regarding the scene. We propose to generate such explanations by identifying the particular question premise(s) that do not apply to an image.

By construction, irrelevant images in the QRPE dataset are picked on the basis of negating a single premise – we now use our dataset to train models to detect false premises, and use the premises classified as irrelevant to generate templated natural language explanations.

Fig. 5 illustrates the task setup for false premise detection. Given a question-image pair, say “*What color is the cat’s tie?*”, the objective is to identify which (if any) question premises are not grounded in the image, in this case both <cat> and <tie>. Alternatively, for the question “*What*

kind of building is the large white building?”, both premises $\langle \textit{building, large} \rangle$ and $\langle \textit{building, white} \rangle$ are true premises grounded in the image.

We train a simple false premise detection model for this task. Our model is a multilayer perceptron that takes one-hot encodings of premises and VGGNet (Simonyan and Zisserman, 2014) image features as input to predict whether the premise is grounded in the image or not. We trained our false premise detection model (FPD) model on all premises in the QRPE dataset.

Our FPD model achieves an accuracy of 61.12% on the QRPE dataset. In Fig. 5, we present qualitative results of our premise classification and explanation pipeline. For the question “*What color is the cat’s tie?*”, the model correctly recognizes ‘cat’ and ‘tie’ as false premises, and we generate statements in natural language indicating the same. Thus, determining question relevance by reasoning about each premise presents the opportunity to generate simple explanations that can provide valuable feedback to the questioner, and help improve model trust.

6 Premise-Based Visual Question Answering Data Augmentation

In this section, we develop a premise-based data augmentation scheme for VQA that generates simple, templated questions based on premises present in complex visually-grounded questions from the VQA (training) dataset.

Using the pipeline presented in Section 3, we extract premises from questions in the VQA dataset and apply a simple templated question generation strategy to transform premises into question and answer pairs. Note that because the source questions come from sighted humans about an image, we do not need to filter out binary or counting questions in order to avoid false premises as in Section 3. We do however filter based on SPICE similarity between the generated and source questions to avoid generating duplicates.

We design templates for each type of premise – first-order (e.g. ‘ $\langle \textit{man} \rangle$ ’ – “*Is there a man?*” *Yes*), second-order (e.g. ‘ $\langle \textit{man, walking} \rangle$ ’ – “*What is the man doing?*” *Walking*, and ‘ $\langle \textit{car, red} \rangle$ ’ – “*What is the color of the car?*” *Red*), and third-order (‘ $\langle \textit{man, holding, racket} \rangle$ ’ – “*What is the man holding?*” *Racket*, “*Who is holding the racket?*” *Man*). This process transforms implicit premise concepts which previously had to

Training Data	Other	Number	Yes	No	Total
Source	123,817	29,698	57217	35842	246,574
Premise	137,483	1,850	387,941	0	527,274

Table 2: Answer type distribution of source and premise questions on the Compositional VQA train set.

be learned as part of understanding more complex questions into simple, explicit training examples that can be directly supervised.

Fig. 6 shows sample premise questions produced from source VQA questions using our pipeline. We note that the distribution of premise questions varies drastically from the source VQA distribution (see Table 2).

We evaluate multiple models with and without premise augmentation on two splits of the VQA dataset - the standard split and the compositional split of Agrawal et al. (2017). The compositional split is specifically designed to test a model’s ability to generalize to unseen/rarely seen combinations of concepts at test time.

Augmentation Strategies. We evaluate the Deeper LSTM model of Lu et al. (2015) on the standard and compositional splits with two augmentation strategies - All which includes the entire set of premise questions and Top-1k-A which includes only questions with answers in the top 1000 most common VQA answers. The results are listed in Table 3. We find minor improvement of 0.34% on the standard split under Top-1k-A premise question augmentation. On the compositional split, we observe a 1.16% gain with Top-1k-A augmentation over no augmentation. In this setting, explicitly reasoning about objects and attributes seen in the questions seems to help the model disentangle objects from their common characteristics.

Other Models. To check the general effectiveness of our approach, we further evaluate Top-1k-A augmentation for three additional VQA models on the compositional split. We find inconsistent improvements for these more advanced models with some improving while others see reductions in accuracy when adding premises.

7 Conclusions and Future Work

In this paper, we made the simple observation that questions about images often contain premises implied by the question and that reasoning about premises can help VQA models respond more in-

<p>What player number is about to swing at the ball?</p> <p>Is there a player number? Yes Is there a ball in the image? Yes Is there a number in the image? Yes</p>	<p>Why is the man looking at the lady?</p> <p>Who is looking at the lady? Man Is there a lady in the image? Yes Is there a man in the image? Yes</p>	<p>How many people are wearing safety jackets?</p> <p>Can you see people in the image? Yes What are the people wearing? Jacket Who is wearing the jacket? People</p>
<p>What is the child sitting on?</p> <p>What is the child doing? Sitting Is there a child in the image? Yes</p>	<p>Where is the pink hat?</p> <p>What is the color of hat? Pink Is there a hat in the image? Yes</p>	<p>What is the item called that the cat is looking at?</p> <p>Is there a cat in the image? Yes Is there an item in the image? Yes</p>

Figure 6: Sample generated premise questions from source questions. Source questions are in bold. Ground-truth answers are extracted using the premise tuples.

	Augmentation	Overall	Other	Number	Yes/No
Standard	None	54.23	40.34	33.27	79.82
	All	53.74	39.28	33.38	79.89
	Top-1k-A	54.47	40.56	33.24	80.19
Comp.	None	46.69	31.92	29.73	70.49
	All	47.63	31.97	30.77	72.52
	Top-1k-A	47.85	32.58	30.59	72.38

Table 3: Accuracy on the standard and compositional VQA validation sets for different augmentation strategies for Deep-LSTM(Antol et al., 2015).

VQA Model	Baseline	+Premises
DeeperLSTM(Lu et al., 2015)	46.69	47.85
HieCoAtt(Lu et al., 2016)	50.17	49.98
NMN(Andreas et al., 2016)	49.05	48.43
MCB(Fukui et al., 2016)	50.13	50.57

Table 4: Overall accuracy of different VQA models on the Compositional VQA test split using Top-1k-A augmentation.

telligently to irrelevant or novel questions.

We develop a system for automatically extracting these question premises. Using these premises, we automatically created a novel dataset for Question Relevance Prediction and Explanation (QRPE) which consists of 53,911 question, relevant image, and irrelevant image triplets. We also train novel question relevance prediction models and show that models that take advantage of premise information outperform models that do not. Furthermore, we demonstrated that questions generated from premises may be an effective data augmentation technique for VQA tasks that require compositional reasoning.

Integrating Question Relevance Prediction and Explanation (QRPE) models with existing VQA systems would form a natural extension to our approach. In this setting, the relevance prediction model would determine the applicability of a ques-

tion to an image, and select an appropriate path of action. If the question is classified as relevant, the VQA model would generate a prediction; otherwise, a question relevance explanation model would provide a natural language sentence indicating which premise(s) are not valid for the image. Such systems would be a step in the direction of making VQA systems move beyond academic settings to real-world environments.

References

- A. Agrawal, A. Kembhavi, D. Batra, and D. Parikh. 2017. C-VQA: A Compositional Split of the Visual Question Answering (VQA) v1.0 Dataset. *arXiv preprint arXiv:1704.08243*.
- 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. Springer.
- Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Dan Klein. 2016. Neural module networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 39–48.
- 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.
- Mohamed Elhoseiny, Scott Cohen, Walter Chang, Brian Price, and Ahmed Elgammal. 2016. Automatic annotation of structured facts in images. *arXiv preprint arXiv:1604.00466*.
- 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. *arXiv preprint arXiv:1606.01847*.
- Lisa Anne Hendricks, Subhashini Venugopalan, Marcus Rohrbach, Raymond Mooney, Saenko Kate, and

- Trevor Darrell. 2016. Deep compositional captioning: Describing novel object categories without paired training data. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Justin Johnson, Ranjay Krishna, Michael Stark, Li-Jia Li, David Shamma, Michael Bernstein, and Li Fei-Fei. 2015. Image retrieval using scene graphs. In *Conference on Computer Vision and Pattern Recognition*, pages 3668–3678.
- Andrej Karpathy and Fei-Fei Li. 2015. Deep visual-semantic alignments for generating image descriptions. In *CVPR*.
- Jin-Hwa Kim, Sang-Woo Lee, Donghyun Kwak, Min-Oh Heo, Jeonghee Kim, Jung-Woo Ha, and Byoung-Tak Zhang. 2016. Multimodal residual learning for visual qa. In *Advances in Neural Information Processing Systems*, pages 361–369.
- Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, Michael Bernstein, and Li Fei-Fei. 2016. *Visual genome: Connecting language and vision using crowdsourced dense image annotations*.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *European Conference on Computer Vision (ECCV)*.
- Jiasen Lu, Xiao Lin, Dhruv Batra, and Devi Parikh. 2015. Deeper lstm and normalized cnn visual question answering model. https://github.com/VT-vision-lab/VQA_LSTM_CNN.
- Jiasen Lu, Caiming Xiong, Devi Parikh, and Richard Socher. 2017. Knowing When to Look: Adaptive Attention via A Visual Sentinel for Image Captioning. *CVPR*.
- Jiasen Lu, Jianwei Yang, Dhruv Batra, and Devi Parikh. 2016. Hierarchical question-image co-attention for visual question answering. In *Advances In Neural Information Processing Systems*, pages 289–297.
- Mateusz Malinowski and Mario Fritz. 2014. A multi-world approach to question answering about real-world scenes based on uncertain input. In *Advances in Neural Information Processing Systems*, pages 1682–1690.
- Mateusz Malinowski, Marcus Rohrbach, and Mario Fritz. 2015. Ask your neurons: A neural-based approach to answering questions about images. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1–9.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119.
- Dong Huk Park, Lisa Anne Hendricks, Zeynep Akata, Bernt Schiele, Trevor Darrell, and Marcus Rohrbach. 2016. Attentive explanations: Justifying decisions and pointing to the evidence. *arXiv preprint arXiv:1612.04757*.
- Arijit Ray, Gordon Christie, Mohit Bansal, Dhruv Batra, and Devi Parikh. 2016. Question relevance in vqa: Identifying non-visual and false-premise questions. In *EMNLP*.
- Olga Russakovsky, Jia Deng, Jonathan Krause, Alex Berg, and Li Fei-Fei. 2012. The ImageNet Large Scale Visual Recognition Challenge 2012 (ILSVRC2012). <http://www.image-net.org/challenges/LSVRC/2012/>.
- Sebastian Schuster, Ranjay Krishna, Angel Chang, Li Fei-Fei, and Christopher D Manning. 2015. Generating semantically precise scene graphs from textual descriptions for improved image retrieval. In *Proceedings of the Fourth Workshop on Vision and Language*, pages 70–80.
- K. Simonyan and A. Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556.
- Peng Wang, Qi Wu, Chunhua Shen, Anton van den Hengel, and Anthony Dick. 2016. Fvqa: Fact-based visual question answering. *arXiv preprint arXiv:1606.05433*.
- Qi Wu, Chunhua Shen, Lingqiao Liu, Anthony Dick, and Anton van den Hengel. 2016. What value do explicit high level concepts have in vision to language problems? In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 203–212.
- Danfei Xu, Yuke Zhu, Christopher Choy, and Li Fei-Fei. 2017. Scene graph generation by iterative message passing. In *Computer Vision and Pattern Recognition (CVPR)*.
- C Lawrence Zitnick, Aishwarya Agrawal, Stanislaw Antol, Margaret Mitchell, Dhruv Batra, and Devi Parikh. 2016. Measuring machine intelligence through visual question answering. *arXiv preprint arXiv:1608.08716*.