# Generating Natural Language Explanations for Visual Question Answering Using Scene Graphs and Visual Attention

**Shalini Ghosh**[*],    **Giedrius Burachas**[†],    **Arijit Ray**[†],    **Avi Ziskind**[†]

[*]`firstname.lastname@samsung.com,`    [†]`firstname.lastname@sri.com`

[*]Samsung Research, [†]SRI International

## Abstract

In this paper, we present a novel approach for the task of eXplainable Question Answering (XQA), i.e., generating natural language (NL) explanations for the Visual Question Answering (VQA) problem. We generate NL explanations comprising of the evidence to support the answer to a question asked to an image using two sources of information: (a) annotations of entities in an image (e.g., object labels, region descriptions, relation phrases) generated from the scene graph of the image, and (b) the attention map generated by a VQA model when answering the question. We show how combining the visual attention map with the NL representation of relevant scene graph entities, carefully selected using a language model, can give reasonable textual explanations without the need of any additional collected data (explanation captions, etc). We run our algorithms on the Visual Genome (VG) dataset and conduct internal user-studies to demonstrate the efficacy of our approach over a strong baseline. We have also released a live web demo showcasing our VQA and textual explanation generation using scene graphs and visual attention.[1]

## 1 Introduction

Visual Question Answering (VQA) [Antol *et al.*, 2015], the task of answering natural language questions on images, has garnered a lot of interest as an AI-complete task. While impressive strides have been made on this task using deep networks [Kazemi and Elqursh, 2017; Teney *et al.*, 2017], they are notorious for being opaque/black-boxed to a non-expert user, thus making it hard to understand when/why it predicts an incorrect answer. There have been attempts to make VQA systems more human-like [Ray *et al.*, 2016] and on how they can hold conversations if one desires further questioning [Ray, 2017] [Das *et al.*, 2016]. However, VQA/conversational agents still cannot explain in natural language why they made a certain decision. This raises issues of trust and reliability since the user cannot judge when to trust the predictions of the model

---

[*]Work done while the first author was at SRI International.
[1]https://xai.nautilus.optiputer.net/

or not. In this paper, we focus on a natural language solution to the eXplainable Question Answering (XQA) task, the task of explaining the answer that was provided — specifically, our goal is to automatically generate a natural language sentence that provides evidence to support the answer predicted by a VQA model. Ideally, such an explanation will help people judge and/or trust the answer provided better.
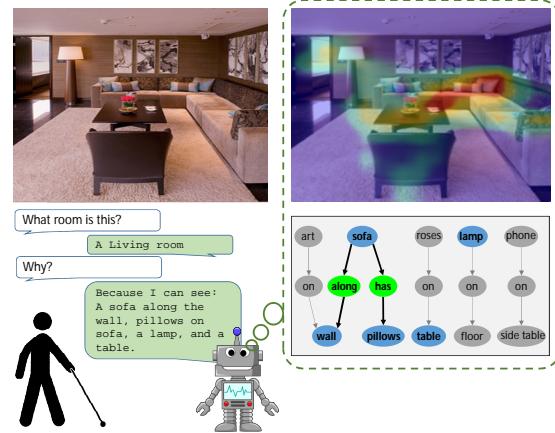


Figure 1: Example Natural Language explanation for an answer to a visually-grounded question. Our approach uses the attention map generated by a visual question-answering model (top right) to identify the relevant components of a scene graph (bottom right) that can be used to provide a justification for the answer given by the model.

Figure 1 illustrates a summary of our objective. Asked a question, "What room is this?", to which the answer is "living room", the visually impaired user requests for an explanation in order to ensure the answer is correct. Our Explainable VQA Agent uses the attention heatmap (regions of interest as suggested by the model while answering the question) to pick out relevant information from an annotated image scene graph to generate a natural language explanation phrase — "Because I can see: sofa along the wall, pillows on sofa, a lamp, and a table." Our proposed algorithm uses the visual attention map as a guide to identify the relevant entities from the scene graph (where entities could be objects, relations or descriptions), and then uses NLP techniques using language models to compose

a NL representation of those entities to generate a NL explanation for the VQA answer. By conducting a small-scale user study, we show evidence that such an approach can generate reasonable textual explanations for answers to questions on images in the Visual Genome Dataset.

Section 2 covers the relevant background and some related work. Section 3 outlines the core algorithms, while Section 4 gives some example explanations generated by our algorithm. We performed some initial experiments using the Visual Genome (VG) dataset to demonstrate the effectiveness of our method — an analysis of those results are described in Section 5. Section 6 discusses related approaches, and finally Section 7 concludes the paper and gives an overview of possible future work in this area.



Figure 2: Salient parts of an Visual Genome image highlighted by the attention layers of our VQA Model, corresponding to the question/answer pair "What is this game? Tennis".

## 2 Background

In this section, we briefly describe some of the key models and concepts used in this paper.

### 2.1 Visual Question Answering and Attention Layers

Visual Question Answering is the task of answering natural language questions about an image. This requires simultaneous understanding of textual and visual semantics — deep neural networks have made impressive strides at this task. We use the VQA architecture outlined in Figure 3. Our model takes as input a 224 x 224 RGB image and a question of at most 15 words. The image is encoded using a ResNet152 [He *et al.*, 2015] to get a 7x7x2048 image feature representation. The question is encoded using an LSTM which takes in the GloVe [Pennington *et al.*, 2014] word embeddings of the words, one word at a time. The final LSTM state is used to represent the question features. The attention layer takes in the question and image feature representations and outputs a set of weights to attend on the image features. The weighted image features, concatenated with the question representation, is used to predict the final answer from a set of 3000 answer choices.

### 2.2 Scene Graph

The scene graph for an image is the graphical representation of its contents, where the nodes are the depicted objects and the edges are the relationships between them (e.g. as shown in Figure 1). Some scene graphs also contain region descriptions (more detailed annotations of an object, or a description of a
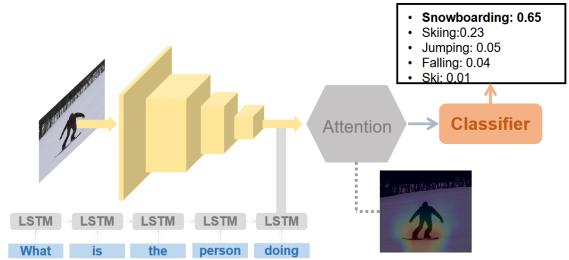


Figure 3: We use a simplistic VQA architecture which is very similar to the SA3 model by Google. The model takes as input a 224 x 224 image and a natural language question (words encoded as 300 dimensional GloVe vectors) and outputs a sigmoidal probability distribution over 3000 answer choices. The attention layer weighs a 7 x 7 x 2048 convolutional map, which is weighted averaged to get a 7 x 7 heatmap. The 7 x 7 heatmap is resized to 224 x 224 to get the attention heatmap for the image.

region containing multiple interacting objects). For example, the Visual Genome data has region descriptions in addition to labeled objects/relations, a sample of which is shown in Figure 4. However, most methods for generating scene graphs [Lu *et al.*, 2016] generate graphs with only objects and relations, but without region descriptions. Depending on the type of scene graph available, we designed two variants of the explanation candidate generation models: (1) one that generates NL explanations based on region descriptions, and (2) one that generates NL explanations based on objects and relations.

### 2.3 Web Language Model

Given a sequence of words, language models estimate the probability of observing another sequence of words following it. There are various language models available, typically trained on large scale corpora (e.g., Web news data, Wikipedia), which give robust estimates of conditional probabilities of observing one text segment in the context of another text segment [Józefowicz *et al.*, 2016]. We use the Web language model service AzureLM[2] to get robust estimates of conditional probabilities of text segments.

## 3 Algorithm

One of the key insights in our approach is the observation that the scene graph of an image can be very useful in generating explanations, something that has not been explored in previous work on generation visual explanations. We use the scene graph to retrieve the set of relevant entities for an image, where an entity refers to either an object, relation or region description, and then consider the natural language phrases corresponding to those entities. Note that each entity in the scene graph has an associated bounding box. For example, for the tennis image in Figure 4, the bounding box shown in the figure is associated with the region description "a tennis player hitting a ball".
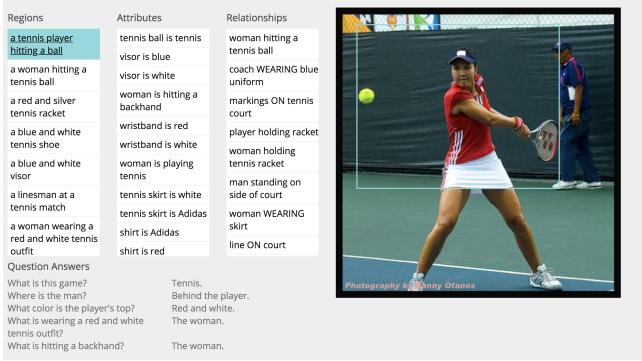
Figure 4: Scene graph of an image from Visual Genome data, showing object attributes, relation phrases and descriptions of regions.

For a given image and a question/answer pair, we use the heatmap generated by the visual attention layer to identify the parts of the image that are relevant and retrieve the bounding boxes with a high degree of overlap with these regions. For example, for the attention map shown in Figure 2, corresponding to the question/answer pair "What is this game? Tennis", the region of the image containing the tennis racket is highlighted, thus identifying this object within the scene graph as relevant for the explanation. We rank the most important entities from the scene graph using a score function, which includes, among other things: (1) The degree of overlap of the bounding box with the active region of the attention map, and (2) An estimate of the relevance of the NL representation of the entity w.r.t. the question/answer.

Our composite score helps us identify entities with NL representations that are deemed relevant to the explanation of the question/answer pair, both from the point of view of the visual attention model as well as the language model. E.g., for Figure 4, a region description with high relevance score (from both the visual model and language model) is "a tennis player hitting a ball" — we use this high-scoring region description to generate the final explanation "The picture shows: a tennis player hitting a ball".
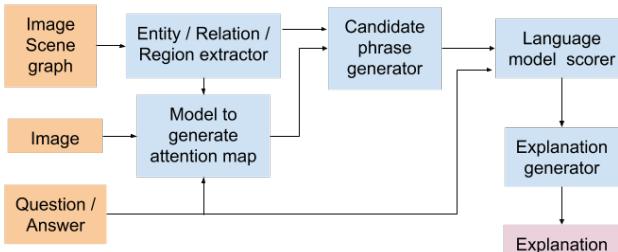


Figure 5: Proposed Explanation Generation Pipeline.

Figure 5 shows the overall flow of explanation generation used in our approach, where we generate explanations based on inferences by a language model and the visual attention layer. When we have region descriptions corresponding to an image, along with the visual attention model, we use the approach in Algorithm 1 to generate the explanations. The score function used to rank regions is:

$$
\begin{aligned}
score(D, QA) \quad = \quad & attentionScore(R(D)|QA) \\
\times \quad & lmScore(D|QA) \\
\times \quad & sqrt(len(D)) \\
\times \quad & 1/log(area(R(D))) \quad (1)
\end{aligned}
$$

where D corresponds to a region description, $R(D)$ is the bounding box of that region, $attention(R(D)|QA)$ is the total attention score from the visual attention layer for region R(D) for question/answer Q/A, $lmScore(D|QA)$ is the language model score of text of D in the context of the text of QA, $len(D)$ is the length of description D and $area(R(D))$ is the area of R(D). Given two regions for related areas, this score function selects the region with tighter bounding box and longer (i.e., richer) description. The slot filler explanation generation in this case simply adds a prefix "The picture shows:" to the generated explanation — in the future, we plan to explore other more complex slot-filter templates (if necessary).

---

**Algorithm 1** Generate XQA explanations using regions

---

**Require:** Image I, question Q, answer A, regions RS in I with descriptions DS, visual (attention) and language (LM) scoring models.
**Ensure:** Explanations E for Q/A, sorted by relevance.
  **for all** Description D in DS, for region R in RS **do**
    Compute relevance of D to Q/A using score(D, QA) defined in Equation 1, using attention and LM models.
  **end for**
  D' = Descriptions D sorted by decreasing score(D, QA)
  **return** E = Slot-filler explanations generated from D'

---

When the scene graph of an image just has objects and relations (and no region descriptions, like in the Visual Genome data), then we use Algorithm 2 to generate the explanations. In this algorithm, we find out the relevant objects/relations and then use a graph traversal algorithm to find a set of connected relations that form a connected component in the graph — we use the connected component to generate a "descriptive" explanation (outlined in Algorithm 3). Thus, when region descriptions are not available, we can instead use these descriptive explanations for the connected component of relations relevant to the question/answer.

## 4 Example Explanations

In this section, we illustrate our algorithm using the Visual Genome (VG) data, providing a few examples in which the multi-modal explanation generation algorithm performs well. For the tennis image and associated Q/A, Algorithm 1, the region-based algorithm gave the following results using the multi-modal explanation generation approach:

**Algorithm 2** Generate XQA explanations using objects/relations

---

**Require:** Image I, question Q, answer A, objects or relations OS in the scene graph, visual (attention) and language (LM) models.
**Ensure:** Explanations E for Q/A, sorted by relevance.
  **for all** Object O in OS **do**
    Compute relevance of O to Q/A using score(O, QA) defined in Equation 1, using attention and LM models.
  **end for**
  O' = Objects O sorted by decreasing score(O, QA)
  **return** E = Descriptive explanations generated from O' using Algorithm 3.

---

**Q/A**: What is this game? Tennis. **Explanation**:
1. The picture shows: a tennis court
2. The picture shows: a tennis player hitting a ball
3. The picture shows: a woman hitting a tennis ball
4. The picture shows: a red and silver tennis racket
5. The picture shows: a blue and white tennis shoe

In this example, explanations 1–4 are relevant. Example 5, while mentioning the relevant concept of a tennis shoe, does not provide as satisfactory an explanation. Consider another image of a set of people at a crosswalk (Figure 6) — for this figure, we get the following results:



Figure 6: VG image of people at crosswalk.

**Q/A**: What is across the street? Other people. **Explanation**:
1. The picture shows: group of people across the street
2. The picture shows: buildings at the end of the street
3. The picture shows: people across street near crosswalk
4. The picture shows: a street lamp
5. The picture shows: people waiting to cross the street

As we can see in this example, explanations 1, 3 and 5 are relevant. The success of these two examples suggests that that getting relevant results in the top 5 explanations.

Our baseline method is to generate a similar explanation but without the attention heatmap guidance from the VQA model. Without the attention heatmap, the sentence generated will still be relevant to the image and question asked, but will not

**Algorithm 3** DFSSortedWithEmit(N): Generate descriptive explanations using objects/relations for graph rooted at N

---

**Require:** Graph G = (OS, RS), OS is set of objects (nodes), RS is set of relations (edges), maximum number of objects used in explanation kNumTermsInExplanation, language (LM) model.
**Ensure:** Explanations E for Q/A.
  LR = list of relations RS in decreasing order of LM model score
  LO = list of objects OS in decreasing order of LM model score
  Explanation list EL = []
  **for all** Relation R(O, O') in LR **do**
    **if** O is in LO and O is not marked **then**
      Phrase P = DFSSortedWithEmit(O)
      Add phrase P to EL
    **end if**
  **end for**
  **if** size(EL) < kNumTermsInExplanation **then**
    **for all** For unmarked object O in LO: **do**
      Phrase P' = DFSSortedWithEmit(O)
      Add phrase P' to EL until size(EL) = kNumTermsInExplanation
    **end for**
  **end if**
  Create explanations E with phrases from list EL using slot-filling
  **return** E

---

be explaining the decision of the model since the explanation generation is not tied to the model decision process in any way. Thus, this acts as a strong baseline since any relevant sentence to the image shouldn't be mistaken for an explanation of why a model predicted a certain answer for a question.

We evaluate how the results for this baseline, i.e., in Equation 1 where the $attentionScore$ is not used. In that case, we expect the results to be of poorer quality. For the tennis example, the following explanations were generated using the attention map:



Figure 7: VG image of crosswalk.

**Q/A**: Question considered: Why is the woman holding a racket? Answer: To hit the ball. **Explanation**:
1. The picture shows: the tennis racket of the player
2. The picture shows: a red and silver tennis racket
3. The picture shows: a woman holding a tennis racket
4. The picture shows: a woman hitting a tennis ball
5. The picture shows: a woman hitting a backhand

As we can see, explanations 4 and 5 are relevant explanations. In comparison, when we don't use the attention map, we typically find that none of the explanations are quite relevant to the answer:

**Q/A**: Why is the woman holding a racket? To hit the ball. **Explanation**:
1. The picture shows: a tennis ball
2. The picture shows: a yellow tennis ball
3. The picture shows: a small tennis ball
4. The picture shows: a red and silver tennis racket
5. The picture shows: the tennis racket of the player

To quantify the importance of the attention map in the multi-modal NL explanation generation algorithm, we performed an A/B test with human raters — the analysis of the ratings, outlined in the next section, validates our assumption that using the multi-modal approach gives us better explanations over using one modality (e.g., linguistic analysis) alone.

We next show some examples of explanations generated by using the story-like algorithm with objects and relations extracted from the scene graph without using region descriptions. In each case, we show the top explanation generated using the story-like explanation generation algorithm, as outlined in Algorithms 2 and 3:



Figure 8: VG image of office.

**Q/A**: Where was this picture taken? At the intersection (Figure 7)
**Explanation**: The picture shows crosswalk on road and in front of man, car parked on road, tree next to road, sign next to road, bike next to car, building with window, walk sign.

**Q/A**: Where was this picture taken? In an office (Figure 8)

**Explanation**: The picture shows keyboard with keys, filing cabinet with drawer, bag on desk, picture on wall, outlet on wall, pen on desk, mouse next to keyboard, filing cabinet with handle, cable on floor, cables on floor.

**Q/A**: Where was this picture taken? In a dining room (Figure 9)
**Explanation**: The picture shows chair with leg, food in bag, liquid in glass, fork on plate, bottle with logo, crumb on plate, bag rests on bowl, bar stool, cover, pan.

## 5 Experiments



Figure 9: VG image of dining room.

We conducted an initial quantitative user-study to demonstrate that human users were satisfied with explanations generated by our multi-modal algorithm. In this evaluation, the ratings were performed internally within our research group, however, in future experiments, we plan to use Amazon Mechanical Turk. Four internal workers rated about 220 questions — for each question, explanations were generated using both the multi-modal algorithm and the NL-only approach (as the baseline). In the latter case, visual attention was not used in the explanation generation. Overall, close to 2K explanations were rated. The authors did not know beforehand which of the explanations came from the baseline algorithm vs. our multi-modal approach, so as to not bias ratings in any way.

We used 3 metrics for evaluation in our small-scale study:

1. Explanation score: Each explanation was rated using a relevance score ranging from -5 to +5, where -5 indicates irrelevant explanation, 0 indicates redundant explanation and +5 corresponds to relevant non-redundant explanation. A negative score between 0 and -5 indicates degree of irrelevance, while a positive score between 0 and +5 indicates degree of relevance.

   Each question/answer pair was also rated, according to the degree of explainability of the question/answer, on score of 1 to 5 — 1 indicates that the question/answer pair is difficult to explain (e.g., a question/answer like "Q: What color is the shirt? A: Red"), while 5 indicates that the question is easy to explain (e.g., a question/answer like "Q: Where was this picture taken? A: On the beach").

| Type | Win | Loss | Tie |
|------|-----|------|-----|
| Explanation score | 52% | 28% | 20% |
| Position score | 55% | 30% | 15% |
| Number score | 54% | 24% | 22% |

Table 1: Statistics of multi-modal approach compared to NL-only approach (Win ⇒ multi-modal wins).

When an algorithm generates explanations for a question/answer pair, we score each explanation using the relevance score, multiply that score by a position weight (so that explanations in higher positions get higher weight), and finally scale that using the question explainability (so that explanations for more explainable questions are given higher score).

2. Position of first relevant explanation: Given a set of generated explanations, this metric compares the position of the first relevant explanation.

3. Number of relevant explanations: Given a set of generated explanations, this metric measures the number of relevant explanations in the top-5 generated explanations.

In all 3 cases, the combined multi-modal algorithm outperformed the baseline NL-only approach. The results are summarized in Table 1: considering explanation score, multimodal was better in ≈ 52% cases; based on position on first relevant score, multi-modal was better in ≈ 55% cases; while in ≈ 54% cases, multi-modal was better than NL-only according to number of relevant explanations. These initial quantitative results demonstrate that humans are more satisfied with our explanations when those explanations came from the multi-modal algorithm that was actually tied to the inference procedure of the visual model.

## 6 Related Work

There is a large body of literature on automatic generation of different types of machine explanations, e.g., explanations for recommendation systems [Costa *et al.*, 2017], affordances from images [Chuang *et al.*, 2017], and robotics [Sridharan *et al.*, 2016]. Explanation generation has also been explored in the areas of planning [Fox *et al.*, 2017], interactive model debugging [Kulesza *et al.*, 2015], autonomous systems [Langley *et al.*, 2017], mobile robotics [Rosenthal *et al.*, 2016], expert systems [Swartout *et al.*, 1991], tactical behavior modeling [van Lent *et al.*, 2004], etc.

In this paper, we focus on NL explanations for the visual question answering task, where previous work has been done by Hendricks et al. [Hendricks *et al.*, 2016] and others. A closely-related work to our approach is VQA-X, a method for generating explanation datasets [Park *et al.*, 2017], where the authors propose a multi-modal methodology for simultaneously generating visual and textual explanations. Another related work for generating NL explanations using a multi-modal approach [Park *et al.*, 2018] qualitatively show cases where visual explanation is more insightful than textual explanation (and vice versa), demonstrating that multi-modal explanation models offer significant benefits over uni-modal approaches. Note that both these approaches rely on getting a large corpus of explanations from human annotators for training the models. In our proposed method, we don't need manually generated explanation data, we use already available annotations from scene graphs only.
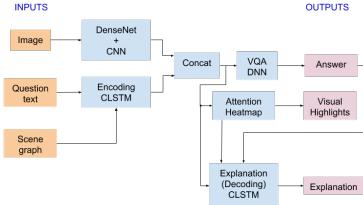


Figure 10: Future Explanation Generation Pipeline.

## 7 Conclusions and Future Work

In this paper, we presented a multi-modal approach for generating natural language explanations for the visual question answering task, using both visual and linguistic modalities, without collecting any additional data. We also showed empirically how the multi-modal approach gives better explanations than using one modality (e.g., linguistic analysis) alone.

We also plan to look into training effective models for learning to generate explanations while predicting answers given an image, context about the image, and a question. Image context may involve scene graphs and we can use Contextual LSTMs (CLSTMs) [Ghosh *et al.*, 2016] to encode this additional information. The main benefit of training a pipeline to do explanation generation (as suggested in Figure 10), instead of using hard-coded algorithms (as suggested in Figure 5), is that models have higher flexibility to learn complicated common-sense semantic information, if that is necessary to generate a satisfactory explanation.

We would also like to explore some other improvements to our system, namely: (1) Infer super-categories of relations and objects (e.g., obtained from hypernyms in wordnet), add them as a preface to the explanation list. (2) Explore the use of embeddings, e.g., skip-thought vectors [Kiros *et al.*, 2015] to find the similarity of entity descriptions (e.g., object attributes, relation phrases or region descriptions) with the salient parts of the question/answer text. (3) Conducting a study to check if our textual explanations help humans predict VQA accuracy for a given image-question pair.

# References

[Antol *et al.*, 2015] Stanislaw Antol, Aishwarya Agrawal, Ji-asen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zit-nick, and Devi Parikh. VQA: Visual Question Answering. In *International Conference on Computer Vision (ICCV)*, 2015.

[Chuang *et al.*, 2017] Ching-Yao Chuang, Jiaman Li, Anto-nio Torralba, and Sanja Fidler. Learning to act properly: Predicting and explaining affordances from images. *CoRR*, abs/1712.07576, 2017.

[Costa *et al.*, 2017] Felipe Costa, Sixun Ouyang, Peter Dolog, and Aonghus Lawlor. Automatic generation of natural language explanations. *CoRR*, abs/1707.01561, 2017.

[Das *et al.*, 2016] Abhishek Das, Satwik Kottur, Khushi Gupta, Avi Singh, Deshraj Yadav, José M. F. Moura, Devi Parikh, and Dhruv Batra. Visual dialog. *CoRR*, abs/1611.08669, 2016.

[Fox *et al.*, 2017] Maria Fox, Derek Long, and Daniele Maga-zzeni. Explainable planning. *CoRR*, abs/1709.10256, 2017.

[Ghosh *et al.*, 2016] Shalini Ghosh, Oriol Vinyals, Brian Strope, Scott Roy, Tom Dean, and Larry P. Heck. Con-textual LSTM (CLSTM) models for large scale NLP tasks. *CoRR*, abs/1602.06291, 2016.

[He *et al.*, 2015] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-nition. *CoRR*, abs/1512.03385, 2015.

[Hendricks *et al.*, 2016] Lisa Anne Hendricks, Zeynep Akata, Marcus Rohrbach, Jeff Donahue, Bernt Schiele, and Trevor Darrell. Generating visual explanations. *CoRR*, abs/1603.08507, 2016.

[Józefowicz *et al.*, 2016] Rafal Józefowicz, Oriol Vinyals, Mike Schuster, Noam Shazeer, and Yonghui Wu. Exploring the limits of language modeling. *CoRR*, abs/1602.02410, 2016.

[Kazemi and Elqursh, 2017] Vahid Kazemi and Ali Elqursh. Show, ask, attend, and answer: A strong baseline for vi-sual question answering. *arXiv preprint arXiv:1704.03162*, 2017.

[Kiros *et al.*, 2015] Ryan Kiros, Yukun Zhu, Ruslan Salakhut-dinov, Richard S. Zemel, Antonio Torralba, Raquel Urtasun, and Sanja Fidler. Skip-thought vectors. In *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 2*, NIPS'15, 2015.

[Kulesza *et al.*, 2015] Todd Kulesza, Margaret Burnett, Weng-Keen Wong, and Simone Stumpf. Principles of ex-planatory debugging to personalize interactive machine learning. In *Proceedings of the 20th International Confer-ence on Intelligent User Interfaces*, 2015.

[Langley *et al.*, 2017] Pat Langley, Ben Meadows, Mohan Sridharan, and Dongkyu Choi. Explainable agency for intelligent autonomous systems. In *AAAI*, 2017.

[Lu *et al.*, 2016] Cewu Lu, Ranjay Krishna, Michael S. Bern-stein, and Fei-Fei Li. Visual relationship detection with language priors. *CoRR*, abs/1608.00187, 2016.

[Park *et al.*, 2017] Dong Huk Park, Lisa Anne Hendricks, Zeynep Akata, Anna Rohrbach, Bernt Schiele, Trevor Dar-rell, and Marcus Rohrbach. Attentive explanations: Jus-tifying decisions and pointing to the evidence (extended abstract). *CoRR*, abs/1711.07373, 2017.

[Park *et al.*, 2018] Dong Huk Park, Lisa Anne Hendricks, Zeynep Akata, Anna Rohrbach, Bernt Schiele, Trevor Dar-rell, and Marcus Rohrbach. Multimodal explanations: Jus-tifying decisions and pointing to the evidence. *CoRR*, abs/1802.08129, 2018.

[Pennington *et al.*, 2014] Jeffrey Pennington, Richard Socher, and Christopher Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 confer-ence on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014.

[Ray *et al.*, 2016] Arijit Ray, Gordon Christie, Mohit Bansal, Dhruv Batra, and Devi Parikh. Question relevance in vqa: identifying non-visual and false-premise questions. *arXiv preprint arXiv:1606.06622*, 2016.

[Ray, 2017] Arijit Ray. *The Art of Deep Connection-Towards Natural and Pragmatic Conversational Agent Interactions*. PhD thesis, Virginia Tech, 2017.

[Rosenthal *et al.*, 2016] Stephanie Rosenthal, Sai P. Selvaraj, and Manuela Veloso. Verbalization: Narration of au-tonomous robot experience. In *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelli-gence (IJCAI)*, 2016.

[Sridharan *et al.*, 2016] Mohan Sridharan, Ben Meadows, and Zenon Colaco. A tale of many explanations: Towards an explanation generation system for robots. In *Proceedings of the 31st Annual ACM Symposium on Applied Computing*, 2016.

[Swartout *et al.*, 1991] W. Swartout, C. Paris, and J. Moore. Explanations in knowledge systems: design for explainable expert systems. *IEEE Expert*, 6(3), 1991.

[Teney *et al.*, 2017] Damien Teney, Peter Anderson, Xi-aodong He, and Anton van den Hengel. Tips and tricks for visual question answering: Learnings from the 2017 challenge. *CoRR*, abs/1708.02711, 2017.

[van Lent *et al.*, 2004] Michael van Lent, William Fisher, and Michael Mancuso. An explainable artificial intelligence system for small-unit tactical behavior. In *Proceedings of the 16th Conference on Innovative Applications of Artifical Intelligence*, IAAI'04, 2004.