

VQA: Visual Question Answering

Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell,
Dhruv Batra, C. Lawrence Zitnick, Devi Parikh

Abstract—We propose the task of *free-form* and *open-ended* Visual Question Answering (VQA). Given an image and a natural language question about the image, the task is to provide an accurate natural language answer. Mirroring many real-world scenarios, such as helping the visually impaired, both the questions and answers are open-ended. Visual questions selectively target different areas of an image, including background details and underlying context. As a result, a system that succeeds at VQA typically needs a more detailed understanding of the image and complex reasoning than a system producing generic image captions. Moreover, VQA is amenable to automatic evaluation, since many open-ended answers contain only a few words or a closed set of answers that can be provided in a multiple-choice format. We provide a dataset containing 100,000's of images and questions and discuss the information it provides. Numerous baselines for VQA are provided and compared with human performance.



1 INTRODUCTION

We are witnessing a renewed excitement in multi-discipline Artificial Intelligence (AI) research problems. In particular, research in image and video captioning that combines Computer Vision (CV), Natural Language Processing (NLP), and Knowledge Representation & Reasoning (KR) has dramatically increased in the past year [14], [7], [10], [31], [22], [20], [42]. Part of this excitement stems from a belief that multi-discipline tasks like image captioning are a step towards solving AI. However, the current state of the art demonstrates that a coarse scene-level understanding of an image paired with word n -gram statistics suffices to generate reasonable image captions, which suggests image captioning may not be as “AI-complete” as desired.

What makes for a compelling “AI-complete” task? We believe that in order to spawn the next generation of AI algorithms, an ideal task should (i) require *multi-modal knowledge* beyond a single sub-domain (such as CV) and (ii) have a well-defined *quantitative evaluation metric* to track progress. For some tasks, such as image captioning, automatic evaluation is still a difficult and open research problem [41], [11], [19].

In this paper, we introduce the task of *open-ended* Visual Question Answering (VQA). A VQA system takes as input an image and a free-form, open-ended, natural-language question about the image and produces an answer as the output. This goal-driven task is applicable to scenarios encountered when visually-impaired users [2] or intelligence analysts actively elicit visual information. Example questions are shown in Fig. 1.

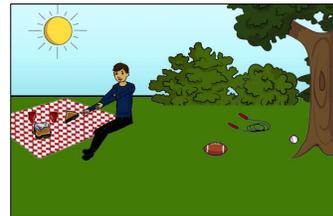
Open-ended questions require a potentially vast set of AI capabilities to answer – fine-grained recognition (e.g., “What kind of cheese is on the pizza?”), object detection (e.g., “How many bikes are there?”), activity recognition (e.g., “Is this man crying?”), knowledge base reasoning (e.g., “Is this a vegetarian pizza?”), and commonsense reasoning (e.g., “Does this person



What color are her eyes?
What is the mustache made of?



How many slices of pizza are there?
Is this a vegetarian pizza?



Is this person expecting company?
What is just under the tree?



Does it appear to be rainy?
Does this person have 20/20 vision?

Fig. 1: Examples of free-form, open-ended questions collected for images via Amazon Mechanical Turk. Note that commonsense knowledge is needed along with a visual understanding of the scene to answer many questions.

have 20/20 vision?”, “Is this person expecting company?”).

VQA [16], [30], [40], [2] is also amenable to automatic quantitative evaluation, making it possible to effectively track progress on this task. While the answer to many questions is simply “yes” or “no”, the process for determining a correct answer is typically far from trivial (for instance in Fig. 1, “Does this person have 20/20 vision?”). Moreover, since questions about images often tend to seek specific information, simple one-to-three word answers are sufficient for many questions. In such scenarios, we can easily evaluate a proposed algorithm by the number of questions it answers correctly. In this paper, we present both an open-ended answering task and a multiple-choice task [37], [28]. Unlike the open-answer task that requires a free-form response, the multiple-choice task only requires an algorithm to pick from a predefined list of possible answers. We anticipate that the multiple-choice task will provide a reasonable *first step* towards VQA, while the open-answer task will become more approachable (and can become more complex) as systems improve.

- S. Antol, A. Agrawal, J. Lu, D. Batra, and D. Parikh are with Virginia Tech.
- M. Mitchell and C. L. Zitnick are with Microsoft Research, Redmond.

We present a large dataset that contains 123,287 images from the Microsoft COCO dataset [27] and a newly created abstract scene dataset [44], [1] that contains 30,000 scenes. The MS COCO dataset has images depicting diverse and complex scenes that are effective at eliciting compelling and diverse questions. We collected a new dataset of “realistic” abstract scenes to enable research focused only on the high-level reasoning required for VQA by removing the need to parse real images. Three questions were collected for each image or scene and each question was answered by three human subjects, along with the subjects’ confidence in their answer. Upon completion, the dataset will contain over 450,000 questions with over 1.3 million answers.¹

While the use of open-ended questions offers many benefits, it is still useful to understand the types of questions that are being asked and which types various algorithms may be good at answering. To this end, we analyze the types of questions asked and the types of answers provided. Through several visualizations, we demonstrate the astonishing diversity of the questions asked. We also explore how the information content of questions and their answers differs from image captions. For baselines, we offer several approaches that use a combination of both text and state-of-the-art visual features [24]. As part of the VQA initiative, we will organize an annual challenge and associated workshop to discuss state-of-the-art methods and best practices.

VQA poses a rich set of challenges, many of which have been viewed as the holy grail of automatic image understanding and AI in general. However, it includes as building blocks several components that the CV, NLP, and KR [4], [6], [26], [29], [3] communities have made significant progress on during the past few decades. VQA provides an attractive balance between pushing the state of the art, while being accessible enough for the communities to start making progress on the task.

2 RELATED WORK

VQA Efforts. Several recent papers have begun to study visual question answering [16], [30], [40], [2]. However, unlike our work, these are fairly restricted (sometimes synthetic) settings with smaller datasets. For instance, [30] only considers questions whose answers come from a predefined closed world of 16 basic colors or 894 object categories. Moreover, many of the questions were synthetically generated via templates. [16] also considers questions generated from templates from a fixed vocabulary of objects, attributes, relationships between objects, *etc.* In contrast, our proposed task involves *open-ended, free-form* questions and answers provided by humans. Our goal is to increase the diversity of knowledge and kinds of reasoning needed to provide correct answers. Critical to achieving success on this more difficult and unconstrained task, our VQA dataset will be *two orders of magnitude* larger than [16], [30] upon completion ($> 150,000$ vs. 2,591 and 1,449 images, respectively). The proposed VQA task has connections to other related work: [40] has studied joint parsing of videos and corresponding text to answer queries on two

datasets containing 15 video clips each. [2] uses crowdsourced workers to answer questions about visual content asked by visually-impaired users.

Text-based Q&A is a well studied problem in the NLP and text processing communities (recent examples being [13], [12], [43], [37]). Other related textual tasks include sentence completion (*e.g.*, [37] with multiple-choice answers). These approaches may provide inspiration for VQA techniques. One key concern in text is the *grounding* of questions. For instance, [43] synthesized textual descriptions and QA-pairs grounded in a simulation of actors and objects in a fixed set of locations. VQA is naturally grounded in images – requiring the understanding of both text (questions) and vision (images). Our questions are generated by humans, making the need for commonsense knowledge and complex reasoning more essential.

Describing Visual Content. Related to VQA are the tasks of image tagging [9], [24], image captioning [25], [15], [33], [7], [14], [42], [10], [20], [31], [22] and video captioning [38], [18], where words or sentences are generated to describe visual content. While these tasks require both visual and semantic knowledge, captions can often be non-specific [42]. The questions posed in VQA require detailed and specific information about the image for which generic image captions are of little use [2].

Other Vision+Language Tasks. Several recent papers have explored tasks at the intersection of vision and language that are easier to evaluate than image captioning, such as coreference resolution [23], [36] or generating referring expressions [21], [35] for a particular object in an image that would allow a human to identify which object is being referred to (*e.g.*, “the one in a red shirt”, “the dog on the left”). While task-driven and concrete, a limited set of visual concepts (*e.g.*, color, location) tend to be captured by referring expressions. As we demonstrate, a richer variety of visual concepts emerge from visual questions and their answers.

3 VQA DATASET COLLECTION

In this section, we describe the Visual Question Answering (VQA) dataset. We begin by describing the real images and abstract scenes used to collect the questions. Next, we describe our process of collecting questions and their corresponding answers. Analysis of the questions and answers gathered as well as baseline results are provided in the following sections.

Real Images. For real images, we use the 123,287 training and validation images from the newly-released Microsoft Common Objects in Context (MS COCO) [27] dataset. The MS COCO dataset was gathered in an attempt to find images containing multiple objects and rich contextual information. Given the visual complexity of these images, they are well-suited for our VQA task. The more diverse our collection of images, the more diverse, comprehensive, and interesting the resultant set of questions and their answers.

Abstract Scenes. The VQA task with real images requires the use of complex and often noisy visual recognizers. To attract researchers interested in exploring the high-level reasoning required for VQA, but not the low-level vision tasks, we create a new abstract scenes dataset [1], [44], [45], [46]. The dataset

1. As of this submission, the dataset currently has 215,150 questions and a total of 430,920 answers for 123,285 MS COCO images and 30,000 questions with 120,810 answers for 10,000 abstract scenes.

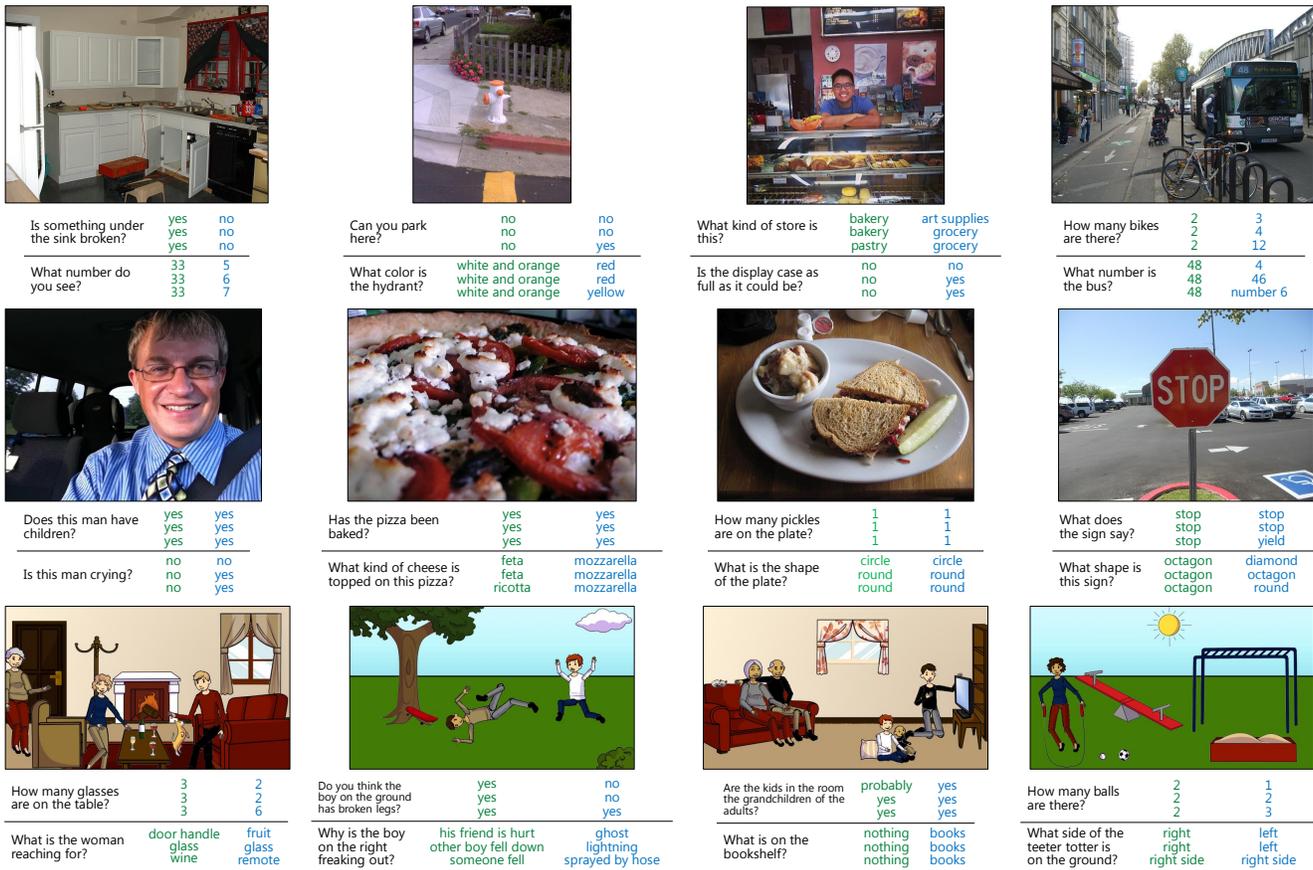


Fig. 2: Examples of questions (black), answers given when looking at the image (green), and answers given when not looking at the image (blue) for numerous representative examples of the dataset. Please see the appendix for more examples.

contains 20 “paperdoll” human models [1] spanning genders, races, and ages with 8 different expressions. The limbs are adjustable to allow for continuous pose variations. The clipart may be used to depict both indoor and outdoor scenes. The set contains over 100 objects and 31 animals in various poses. The use of this clipart enables the creation of more realistic scenes (see bottom row of Fig. 2) that more closely mirror real images than previous papers [44], [45], [46]. Upon completion, 30,000 abstract scenes will be collected. Please see the appendix for the user interface, additional details, and examples.

Captions. The MS COCO dataset [27], [5] already contains five single-sentence captions for all images. Upon completion, our VQA dataset will contain captions for the abstract scenes using the same user interface² for collection.

Questions. Collecting interesting, diverse, and well-posed questions is a significant challenge. Many simple questions may only require low-level computer vision knowledge, such as “What color is the cat?” or “How many chairs are present in the scene?”. However, we also want questions that require commonsense knowledge about the scene, such as “What sound does the pictured animal make?”. Importantly, questions should also *require* the image to correctly answer and not be answerable using just commonsense information, *e.g.*, in Fig. 1, “What is the mustache made of?”. By having a wide variety of question types and difficulty, we may be able

to measure the continual progress of both visual understanding and commonsense reasoning.

We tested and evaluated a number of user interfaces for collecting such “interesting” questions. Specifically, we ran pilot studies asking human subjects to ask questions about a given image that they believe a “toddler”, “alien”, or “smart robot” would have trouble answering. We found the “smart robot” interface to elicit the most interesting and diverse questions. As shown in the appendix, our final interface stated “We have built a smart robot. It understands a lot about images. It can recognize and name all the objects, it knows where the objects are, it can recognize the scene (*e.g.*, kitchen, beach), people’s expressions and poses, and properties of objects (*e.g.*, color of objects, their texture). Your task is to stump this smart robot!”. To bias against generic image-independent questions, subjects were instructed to ask questions that *require* the image to answer.

The same user interface was used for both the real images and abstract scenes. In total, three questions by unique workers were gathered for each image or scene. When writing a question, the subjects were shown the previous questions already asked for that image to help increase the question diversity. In total, the dataset will contain over 450,000 questions when completed.

Answers. Open-ended questions result in a diverse set of possible answers. For many questions, a simple “yes” or “no” response is sufficient. However, other questions may require a

2. <https://github.com/tylin/coco-ui>

short phrase. Multiple different answers may also be correct. For instance, the answers “white”, “tan”, or “off-white” may all be correct answers to the same question. Human subjects may also disagree on the “correct” answer, *e.g.*, some saying “yes” while others say “no”. To handle these discrepancies, we gather three answers for each question from unique workers, while also ensuring that the worker answering a question did not ask it. We ask the subjects to provide answers that are “a brief phrase and not a complete sentence. Respond matter-of-factly and avoid using conversational language or inserting your opinion.” In addition to answering the questions, the subjects were asked “Do you think you were able to answer the question correctly?” and given the choices of “no”, “maybe”, and “yes”. Please see Section 4 for an analysis of the answers provided.

For testing, we offer two modalities for answering the questions: (i) open-answer and (ii) multiple choice.

For the open-answer task, the generated answers are scored using the percentage of the human subjects’ answers that exactly correspond to the generated answer. Before comparison, all responses are made lowercase, numbers converted to digits, and punctuation and articles removed. We avoid using soft metrics such as Word2Vec [32], since they often group similar types of words together that we wish to distinguish, such as “left” and “right”.

For multiple choice, 18 candidate answers are created for each question. As with the open-answer task, accuracy is computed based on the percentage of human subjects that provided an answer that exactly corresponds to the picked answer. We generate a candidate set of correct and incorrect answers from four sets of answers: **Correct**: The set of correct answers. **Plausible**: To generate incorrect, but still plausible, answers we ask three subjects to answer the questions without seeing the image. If three unique answers are not found, we gather additional answers from nearest neighbor questions using a bag-of-words model. The use of these answers helps ensure the image, and not just commonsense knowledge, is necessary to answer the question. **Popular**: These are the 10 most popular answers: “yes”, “no”, “2”, “1”, “3”, “white”, “4”, “red”, “blue”, “black”. The inclusion of the most popular answers makes it more difficult for algorithms to infer the type of question from the set of answers provided, *i.e.*, learning that it is a “yes or no” question just because “yes” and “no” are present in the answers. **Random**: Random answers from other questions. To generate a total of 18 candidate answers, we first find the union of the correct, plausible, and popular answers. The remaining answers are selected from the random set until 18 unique answers are found. The order of the answers is randomized. Example multiple-choice questions are provided in the appendix.

4 VQA DATASET ANALYSIS

In this section, we provide analysis on the questions and answers in the VQA dataset. To gain an understanding of the types of questions asked and answers provided, we visualize the distribution of question types and answers. We also explore how often the questions may be answered without the image using just commonsense information. Finally, we analyze

whether the information contained in an image caption is sufficient to answer the questions.

We emphasize that the creation of a dataset of this scale and richness is a time consuming process, taking months to complete. As of this submission, the VQA dataset has 123,285 images, 215,150 questions, and a total of 430,920 answers (including answers provided by workers while looking at the image and without looking at the image; recall that the latter become options in the multiple-choice question format) for real images from the MS COCO dataset [27]. 30,000 questions with 120,810 answers (again, multiple-choice ones included) for 10,000 abstract scenes have been collected.

4.1 Questions

Types of Question. Given the structure of questions generated in the English language, we can cluster questions into different types based on the words that start the question. Fig. 3 shows the distribution of questions based on the first four words of the questions for both the real images (left) and abstract scenes (right). Interestingly, the distribution of questions is quite similar for both real images and abstract scenes. This helps demonstrate that the type of questions elicited by the abstract scenes is similar to those elicited by the real images. There exists a surprising variety of question types, including “What is...”, “Is there...”, “How many...”, and “Does the...”. Quantitatively, the percentage of questions for different types is shown in Table 3. Several example questions and answers are shown in Fig. 2.

A particularly interesting type of question is the “What is...” questions, since they have a diverse set of possible answers. Please see the appendix for more visualizations showing both the diversity in questions and answers for “What is...” questions.

Lengths. Fig. 4 shows the distribution of question lengths. We can see that most questions range from four to ten words.

4.2 Answers

Typical Answers. Fig. 5 (top) shows the distribution of answers for several question types. We can see that a number of question types, such as “Is the...”, “Are...”, and “Does...” are typically answered using “yes” and “no” as answers. Other questions such as “What is...” and “What type...” have a rich diversity of responses. Other question types such as “What color...” or “Which...” have more specialized responses, such as colors, or “left” and “right”. Please see the appendix for a list of the most popular answers.

Lengths. Most answers consist of a single word, with 86.88%, 8.38%, and 3.25% of answers containing one, two, or three words, respectively. The brevity of answers is not surprising, since the questions tend to elicit specific information from the images. This is in contrast with image captions that generically describe the entire image and hence tend to be longer. The brevity of our answers makes automatic evaluation feasible. While it may be tempting to believe the brevity of the answers makes the problem easier, recall that they are human-provided open-ended answers to open-ended questions. The questions typically require complex reasoning to arrive at these

	Open-Answer			Multiple-Choice*			Yes/No
	100	K 500	1000	100	K 500	1000	K 2
Question	26.91	26.68	26.64	29.60	30.18	30.18	60.11
Image	16.29	16.29	16.29	15.78	15.76	15.78	53.79
Caption	16.08	16.03	16.03	16.11	15.98	15.90	53.22
Q+I	27.28	27.13	26.78	29.70	30.47	30.15	60.02
Q+C	30.34	29.70	29.36	32.29	32.86	32.56	59.16
Q+I+C	30.41	28.91	29.39	32.33	32.24	32.45	59.28

TABLE 2: Results on the open-answer, multiple-choice, and “yes/no” question tasks for various baselines on real images. For comparison, human accuracies on the open-answer task are 25.46 for Question, 61.58 for Q+I, and 38.28 for Q+C; and on the Yes/No task are 83.88, see Table 1. All results are the percentage of answers in agreement with human subjects. See text for details. *The set of questions differ between Open-Answer and Multiple-Choice, so Multiple-Choice accuracies may not be strictly higher than Open-Answer.

Question Type	Open-Answer				
	K = 500			Human	
	Q	Q + I	Q + C	Q	Q + I
what is (14.66)	06.35	07.05	13.77	08.30	47.11
what color (07.80)	16.26	17.05	21.58	19.71	65.65
what kind (02.73)	04.22	05.02	10.00	09.98	39.98
what are (02.48)	08.13	08.59	15.05	10.75	45.36
what type (02.25)	07.41	08.50	14.02	08.80	40.62
is the (09.26)	52.33	52.27	51.10	51.03	83.37
is this (06.67)	50.42	51.49	51.43	48.96	81.82
how many (11.62)	27.21	27.25	29.96	17.73	64.43
are (06.39)	54.80	55.01	53.58	52.59	83.39
does (03.12)	59.53	58.79	59.25	57.13	82.78
where (02.52)	02.10	02.23	04.86	04.73	20.82
is there (02.75)	72.34	72.55	71.57	60.55	84.36
why (01.51)	03.17	03.42	03.67	04.52	08.21
which (01.47)	17.41	17.86	17.94	15.71	39.76
do (01.35)	58.55	58.33	56.94	54.70	80.48
what does (01.28)	07.19	07.19	08.20	08.29	50.26
what time (00.87)	05.56	05.56	05.31	02.67	30.84
who (00.82)	08.02	08.02	07.02	05.66	30.89
what sport (00.78)	40.61	55.99	72.17	13.63	80.18
what animal (00.61)	15.91	17.61	36.17	11.41	75.83
what brand (00.43)	05.01	05.76	10.28	20.41	59.68

TABLE 3: Results on the open-answer task for various question types on real images. Questions types are determined by the one or two words that start the question. The percentage of questions for each type is shown in parentheses. All results are the percentage of answers in agreement with human subjects. See text for details.

when subjects are shown the actual image. This demonstrates that in order to answer the questions correctly, deeper image understanding beyond what image captions typically capture is necessary. In fact, we find that the distributions of nouns, verbs, and adjectives mentioned in captions describing the real images is statistically significantly different from those mentioned in our questions + answers (Kolmogorov-Smirnov test, $p < .001$). See appendix for further details. This motivates the VQA task as a way to learn further information about visual scenes.

5 VQA BASELINES

In this section, we explore the difficulty of the VQA dataset using several different baselines. For reference, if we randomly choose an answer from the training dataset the accuracy is $< 0.5\%$ and if the most popular answer is always selected (“yes”) the accuracy is 16.29%.

For testing, we split our real image dataset into a training and testing dataset. Our training dataset contains 30,000 images and our testing dataset contains 20,000 images.³ For our baselines, we choose the top $K = \{100, 500, 1000\}$ most frequent answers as possible outputs. These sets of answers cover 63%, 75%, and 80% of the test set answers for $K = \{100, 500, 1000\}$, respectively. All baselines train a softmax neural network classifier with a single hidden layer containing 50 units. Six baselines using different inputs are used: **Question:** The top 1,000 words in the questions are used to create a bag-of-words representation. Since there is a strong correlation between the words that start a question and the answer (see Fig. 3), we find the top 10 first, second, and third words of the questions and create a 30 dimensional bag-of-words representation. These features are concatenated to get a 1,030 dimensional input representation. **Caption:** Similar to Table 1, we assume that a human-generated caption is given as input. We use a bag-of-words representation containing the 1,000 most popular words in the captions as the input feature. **Image:** We use the final softmax predictions from AlexNet [24] as our 1,000 dimensional feature. We found this provides better results than the intermediate layers (f_{c6} and f_{c7}). We also report the learned baseline results on **Question+Image (Q+I)**, **Question+Caption (Q+C)**, and **Question+Image+Caption (Q+I+C)** by simply concatenating the features.

For testing, we report the result on two different tasks: Open-answer selects the answer with highest activation from all possible K answers and multiple-choice picks the answer that has the highest activation from the potential answers. As shown in Table 2, the accuracy using only the question is $\sim 27\%$, which demonstrates that the type of question is informative of the answer. As a sanity check, if we select the most popular answer for each question type, the accuracy is 23.08%. We decide on different question types based on first few words of questions and ensure that each question type has at least 125 questions in the dataset. Note that this performs worse than our learned baseline that uses only the question. Table 2 also shows that the results on multiple-choice are generally better than open-answer, as expected. In comparison to human performance, the baseline results are still significantly worse.

In Table 2, the results of (Q + I) are very similar to the results of only using the questions, while the addition of captions (Q + C) provides minimal improvement $\sim 2 - 3\%$. To gain insight into these results, we computed accuracies by question type in Table 3. Interestingly, for question types that require more reasoning, such as “Is the . . .” or “Which . . .”, the scene-level image features do not provide any additional information. However, for questions that can be answered using scene-level information, such as “What sport . . .” or “What animal . . .”, we do see an improvement. Similarly, for questions whose answer may be contained in a generic caption we see improvement, such as “What color . . .” or “What animal . . .”. For all question types, the results are worse than humans.

We also tested our baselines on only binary $K = 2$ (*i.e.*,

³ As of this submission, only 50,000 of the real images had their questions answered.

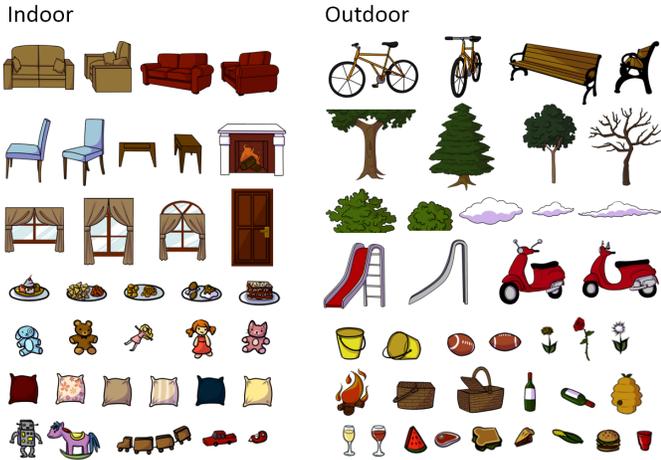


Fig. 13: A small subset of the objects present in the abstract scene dataset.

Help Us Create Clipart Illustrations! (Living/Dining Room)



Fig. 14: The AMT interface for collecting abstract scenes. The light green circles indicate where users can select to manipulate a person’s pose. Different objects may be added to the scene using the folders to the right.

APPENDIX V: USER INTERFACES

In Fig. 15, we show the AMT interface that we used to collect questions for images. Note that we tell them that the robot already knows the answer to the previously asked question(s), inspiring the workers to ask different kinds of questions, thereby increasing the diversity of our dataset.

Fig. 16 shows the AMT interface used for collecting answers to the previously collected questions when subjects were shown the corresponding images. Fig. 17 shows the interface that was used to collect answers to questions when subjects were not shown the corresponding images (*i.e.*, to help in gathering incorrect, but plausible, answers for the multiple-choice task and to assess how accurately the questions can be answered using common sense knowledge alone).

Stump a smart robot! Ask a question about this scene that a human can answer, but a smart robot probably can't!

Updated instructions: Please read carefully

Hide Show

We have built a smart robot. It understands a lot about scenes. It can recognize and name all the objects, it knows where the objects are, it can recognize the scene type (e.g. kitchen, beach), people's expressions and poses, and properties of objects (e.g. the color of objects, their texture). Your task is to stump this smart robot! In particular, it already knows answers to some questions about this scene. We will tell you what those questions are.

Ask a question about this scene that this SMART robot probably can't answer, but any human can easily answer while looking at the scene in the image. IMPORTANT: The question should be about this scene. That is, the human should need the image to be able to answer the question – the human should not be able to answer the question without looking at the image.



Your work will get rejected if you do not follow the instructions below:

- Do not ask questions that are similar to the ones listed below each image. As mentioned, the robot already knows the answers to those questions for the scene in this image. Please ask about something different.
- Do not repeat questions. Do not ask the same questions or the same questions with minor variations over and over again across images. Think of a new question each time specific to the scene in each image.
- Each question should be a single question. Do not ask questions that have multiple parts or multiple sub-questions in them.
- Do not ask generic questions that can be asked of many other scenes. Ask questions specific to the scene in each image.

Below is a list of questions the smart robot can already answer. Please ask a different question about this scene that a human can answer "if" looking at the scene in the image (and not otherwise), but would stump this smart robot:

Q1: What is unusual about this mustache? (The robot already knows the answer to this question.)
 Q2: What is her facial expression? (The robot already knows the answer to this question.)
 Q3: Write your question, different from the questions above, here to stump this smart robot.

prev next Page 2/3

Fig. 15: Our AMT interface for collecting the third question for an image, when subjects were shown previous questions that were collected and were asked to ask a question different from previous questions.

Help Us Answer Questions About Images!

Updated instructions: Please read carefully

Hide Show

Please answer some questions about images with brief answers. Your answers should be how most other people would answer the questions. If the question doesn't make sense, please try your best to answer it and indicate via the buttons that you are unsure of your response.

If you don't follow the following instructions, your work will be rejected.



- Your work will get rejected if you do not follow the instructions below:
- Answer the question based on what is going on in the scene depicted in the image.
 - Your answer should be a brief phrase (not a complete sentence).
 - It is a kitchen? -> "kitchen"
 - For yes/no questions, please just say yes/no.
 - "You bet it is!" -> "yes"
 - For numerical answers, please use digits.
 - "Ten" -> "10"
 - If you need to speculate (e.g. "What just happened?"), provide an answer that most people would agree on.
 - If you don't know the answer (e.g. specific dog breed), provide your best guess.
 - Respond matter-of-factly and avoid using conversational language or inserting your opinion.

Please answer the question using as few words as possible:

Q1: What is unusual about this mustache?
 A1: Write your answer here.

Do you think you were able to answer the question correctly?
 (Clicking an option will take you to the next question.)

no maybe yes Page 1/2

Fig. 16: The AMT interface used to collect answers to a question when subjects were shown the image while answering the question.

Help Us Answer Questions!

Updated instructions: please read carefully

We will show you a series of questions about possibly different scenes. Your task is to answer them. Here's the catch: we will not show you the scene! So how can you answer the question correctly? If you can't, but your job is to provide a plausible answer to the question. What this means is the following: If we show the question alongside your answer to someone else (who also can't see the scene), they should think your answer "could be" correct. Please keep your answer brief! If the question doesn't make sense, please try your best to answer it and indicate via the buttons that you are unsure of your response.

If you don't follow the following instructions, your work will be rejected.

- Instructions:
- Your answer should be a brief phrase (not a complete sentence).
 - It is a kitchen? -> "kitchen"
 - For yes/no questions, please just say yes/no.
 - "You bet it is!" -> "yes"
 - For numerical answers, please use digits.
 - "Ten" -> "10"
 - Respond matter-of-factly and avoid using conversational language.

Please provide a plausible answer to the question using as few words as possible:

Q1: What is unusual about this mustache?
 A1: Write your answer here.

How likely do you think it is that someone else would answer this question with the same answer as yours?
 (Clicking an option will take you to the next question.)

very unlikely somewhat unlikely somewhat likely very likely Page 1/2

Fig. 17: The AMT interface used to collect answers to a question when subjects were not shown the image while answering the question using only commonsense to collect the plausible, but incorrect, multiple-choice answers.

APPENDIX VI: ANSWER DISTRIBUTION

The top 250 answers in our dataset along with their counts and percentage counts are given below. The answers have been presented in different colors to show the different Part-of-Speech (POS) tagging of the answers with the following color code: **yes/no**, **noun**, **verb**, **adjective**, **adverb**, and **numeral**.

“yes” (49027, 18.14%), “no” (38368, 14.2%), “2” (9264, 3.43%), “1” (5095, 1.89%), “3” (4962, 1.84%), “white” (3428, 1.27%), “4” (3378, 1.25%), “red” (3260, 1.21%), “blue” (2627, 0.97%), “black” (2010, 0.74%), “green” (1942, 0.72%), “5” (1803, 0.67%), “yellow” (1593, 0.59%), “brown” (1393, 0.52%), “6” (1372, 0.51%), “0” (1330, 0.49%), “tennis” (1127, 0.42%), “right” (1080, 0.4%), “orange” (1028, 0.38%), “left” (1018, 0.38%), “frisbee” (954, 0.35%), “baseball” (897, 0.33%), “bathroom” (858, 0.32%), “7” (837, 0.31%), “pink” (832, 0.31%), “wood” (801, 0.3%), “pizza” (792, 0.29%), “8” (780, 0.29%), “cat” (738, 0.27%), “water” (605, 0.22%), “dog” (600, 0.22%), “surfing” (590, 0.22%), “skiing” (563, 0.21%), “wii” (554, 0.21%), “elephant” (513, 0.19%), “kite” (506, 0.19%), “skateboarding” (505, 0.19%), “giraffe” (500, 0.19%), “10” (492, 0.18%), “man” (466, 0.17%), “kitchen” (465, 0.17%), “skateboard” (464, 0.17%), “9” (459, 0.17%), “snow” (447, 0.17%), “apple” (433, 0.16%), “stop” (428, 0.16%), “sheep” (422, 0.16%), “sunny” (415, 0.15%), “banana” (415, 0.15%), “umbrella” (406, 0.15%), “grass” (405, 0.15%), “horse” (392, 0.15%), “winter” (391, 0.14%), “12” (382, 0.14%), “bear” (359, 0.13%), “male” (357, 0.13%), “zebra” (349, 0.13%), “bananas” (343, 0.13%), “eating” (340, 0.13%), “female” (338, 0.13%), “purple” (337, 0.12%), “grey” (334, 0.12%), “tan” (331, 0.12%), “black and white” (330, 0.12%), “train” (328, 0.12%), “phone” (317, 0.12%), “wine” (315, 0.12%), “broccoli” (315, 0.12%), “gray” (307, 0.11%), “unknown” (300, 0.11%), “woman” (291, 0.11%), “cake” (289, 0.11%), “cow” (287, 0.11%), “laptop” (280, 0.1%), “beach” (271, 0.1%), “hat” (270, 0.1%), “silver” (269, 0.1%), “night” (255, 0.09%), “motorcycle” (254, 0.09%), “surfboard” (250, 0.09%), “kites” (250, 0.09%), “bird” (237, 0.09%), “maybe” (237, 0.09%), “soccer” (236, 0.09%), “bench” (236, 0.09%), “many” (234, 0.09%), “flowers” (227, 0.08%), “snowboarding” (225, 0.08%), “cell phone” (221, 0.08%), “bus” (217, 0.08%), “nothing” (216, 0.08%), “11” (215, 0.08%), “15” (215, 0.08%), “food” (214, 0.08%), “brick” (213, 0.08%), “20” (211, 0.08%), “christmas” (210, 0.08%), “day” (210, 0.08%), “up” (209, 0.08%), “outside” (206, 0.08%), “down” (205, 0.08%), “teddy bear” (203, 0.08%), “camera” (202, 0.07%), “open” (200, 0.07%), “helmet” (196, 0.07%), “sitting” (186, 0.07%), “fork” (182, 0.07%), “standing” (179, 0.07%), “zoo” (178, 0.07%), “13” (177, 0.07%), “plane” (177, 0.07%), “overcast” (176, 0.07%), “ball” (175, 0.06%), “cows” (175, 0.06%), “tree” (174, 0.06%), “red and white” (173, 0.06%), “cold” (169, 0.06%), “fence” (168, 0.06%), “suitcase” (167, 0.06%), “tile” (164, 0.06%), “glass” (164, 0.06%), “wedding” (163, 0.06%), “bed” (163, 0.06%), “14” (163, 0.06%), “happy” (162, 0.06%), “metal” (162, 0.06%), “lot” (161, 0.06%), “fall” (158, 0.06%), “beer” (158, 0.06%), “girl” (155, 0.06%), “boat” (154, 0.06%), “nike” (152, 0.06%), “bike” (152, 0.06%), “old” (152, 0.06%), “table” (151, 0.06%), “cheese” (149, 0.06%), “asian” (148, 0.05%), “coffee” (147,

0.05%), “flying kite” (145, 0.05%), “sleeping” (143, 0.05%), “horses” (143, 0.05%), “trees” (143, 0.05%), “cloudy” (142, 0.05%), “scissors” (141, 0.05%), “hot dog” (140, 0.05%), “blue and white” (140, 0.05%), “16” (140, 0.05%), “fire hydrant” (139, 0.05%), “chinese” (138, 0.05%), “airport” (137, 0.05%), “fish” (137, 0.05%), “beige” (136, 0.05%), “fruit” (133, 0.05%), “birds” (132, 0.05%), “airplane” (132, 0.05%), “on” (132, 0.05%), “glasses” (131, 0.05%), “microwave” (131, 0.05%), “clock” (130, 0.05%), “sand” (129, 0.05%), “africa” (129, 0.05%), “mirror” (129, 0.05%), “plastic” (128, 0.05%), “summer” (128, 0.05%), “tie” (128, 0.05%), “50” (127, 0.05%), “chocolate” (127, 0.05%), “birthday” (127, 0.05%), “carrots” (126, 0.05%), “truck” (126, 0.05%), “living room” (124, 0.05%), “car” (123, 0.05%), “25” (123, 0.05%), “boy” (122, 0.05%), “walking” (122, 0.05%), “palm” (122, 0.05%), “donuts” (121, 0.04%), “snowboard” (120, 0.04%), “toilet” (120, 0.04%), “donut” (119, 0.04%), “child” (119, 0.04%), “30” (119, 0.04%), “knife” (118, 0.04%), “breakfast” (116, 0.04%), “chair” (116, 0.04%), “round” (116, 0.04%), “mountains” (116, 0.04%), “clear” (116, 0.04%), “bat” (116, 0.04%), “evening” (115, 0.04%), “stripes” (115, 0.04%), “paper” (114, 0.04%), “morning” (114, 0.04%), “carrot” (114, 0.04%), “blonde” (113, 0.04%), “person” (113, 0.04%), “elephants” (112, 0.04%), “watch” (111, 0.04%), “sandwich” (110, 0.04%), “toothbrush” (109, 0.04%), “closed” (109, 0.04%), “afternoon” (108, 0.04%), “usa” (108, 0.04%), “tv” (108, 0.04%), “bicycle” (107, 0.04%), “books” (106, 0.04%), “leaves” (106, 0.04%), “heart” (104, 0.04%), “i dont know” (104, 0.04%), “chicken” (104, 0.04%), “surf board” (103, 0.04%), “hotdog” (103, 0.04%), “daytime” (100, 0.04%), “computer” (100, 0.04%), “tennis racket” (99, 0.04%), “gold” (99, 0.04%), “circle” (99, 0.04%), “dirt” (98, 0.04%), “oranges” (98, 0.04%), “square” (98, 0.04%), “mountain” (98, 0.04%), “ocean” (97, 0.04%), “wall” (97, 0.04%), “cooking” (96, 0.04%), “mouse” (96, 0.04%), “rock” (96, 0.04%), “remote” (95, 0.04%), “100” (95, 0.04%), “playing” (95, 0.04%), “spoon” (95, 0.04%), “18” (94, 0.03%), “bread” (94, 0.03%), “brown and white” (93, 0.03%), “very” (93, 0.03%), “bedroom” (93, 0.03%), “plaid” (93, 0.03%), “window” (93, 0.03%), “hay” (92, 0.03%), “i don t know” (92, 0.03%), “small” (91, 0.03%), “people” (91, 0.03%), “milk” (91, 0.03%), “apples” (90, 0.03%), “towel” (90, 0.03%), “middle” (88, 0.03%), “luggage” (87, 0.03%), “backpack” (87, 0.03%), “skis” (86, 0.03%), “don t know” (85, 0.03%), “ketchup” (85, 0.03%), “giraffes” (84, 0.03%), “vase” (84, 0.03%), “sun” (83, 0.03%), “17” (83, 0.03%), “brushing teeth” (83, 0.03%), “both” (83, 0.03%), “restaurant” (81, 0.03%), “roses” (81, 0.03%), “ground” (81, 0.03%).

APPENDIX VII: ADDITIONAL EXAMPLES

To provide insight into the dataset, we provide additional examples. In Fig. 18, Fig. 19, and Fig. 20, we show a random selection of the VQA dataset for the MS COCO [27] images, abstract scenes, and multiple-choice questions, respectively.

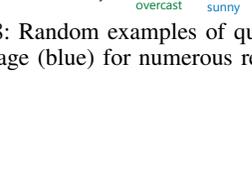
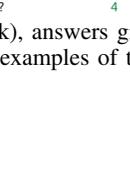
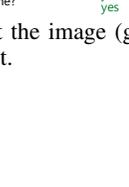
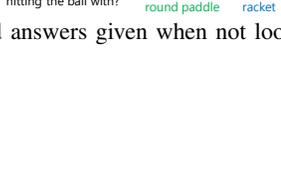
	<p>What part of the body are these worn around?</p> <table border="1"><tr><td>neck</td><td>legs</td></tr><tr><td>neck</td><td>neck</td></tr><tr><td>neck</td><td>wrist</td></tr></table>	neck	legs	neck	neck	neck	wrist		<p>Is it raining?</p> <table border="1"><tr><td>no</td><td>yes</td></tr><tr><td>no</td><td>yes</td></tr><tr><td>no</td><td>Yes</td></tr></table>	no	yes	no	yes	no	Yes		<p>What color are the shoe laces?</p> <table border="1"><tr><td>Blue</td><td>Black</td></tr><tr><td>Blue</td><td>Red</td></tr><tr><td>Light blue</td><td>White</td></tr></table>	Blue	Black	Blue	Red	Light blue	White		<p>How many boats are visible?</p> <table border="1"><tr><td>3</td><td>2</td></tr><tr><td>4</td><td>3</td></tr><tr><td>5</td><td>4</td></tr></table>	3	2	4	3	5	4						
neck	legs																																				
neck	neck																																				
neck	wrist																																				
no	yes																																				
no	yes																																				
no	Yes																																				
Blue	Black																																				
Blue	Red																																				
Light blue	White																																				
3	2																																				
4	3																																				
5	4																																				
	<p>How many ties are there?</p> <table border="1"><tr><td>2</td><td>3</td></tr><tr><td>10</td><td>4</td></tr><tr><td>many</td><td>4</td></tr></table>	2	3	10	4	many	4		<p>Are there any people sitting on the bench?</p> <table border="1"><tr><td>No</td><td>No</td></tr><tr><td>No</td><td>No</td></tr><tr><td>No</td><td>Yes</td></tr></table>	No	No	No	No	No	Yes		<p>What is he sitting on?</p> <table border="1"><tr><td>Skateboard</td><td>Bench</td></tr><tr><td>Skateboard</td><td>Chair</td></tr><tr><td>Skateboard</td><td>Chair</td></tr></table>	Skateboard	Bench	Skateboard	Chair	Skateboard	Chair		<p>How many umbrellas are in the image?</p> <table border="1"><tr><td>4</td><td>2</td></tr><tr><td>4</td><td>2</td></tr><tr><td>4</td><td>9734</td></tr></table>	4	2	4	2	4	9734						
2	3																																				
10	4																																				
many	4																																				
No	No																																				
No	No																																				
No	Yes																																				
Skateboard	Bench																																				
Skateboard	Chair																																				
Skateboard	Chair																																				
4	2																																				
4	2																																				
4	9734																																				
	<p>Does this look like a group of nerds?</p> <table border="1"><tr><td>no</td><td>no</td><td>no</td></tr><tr><td>no</td><td>no</td><td>yes</td></tr><tr><td>no</td><td>yes</td><td>yes</td></tr></table>	no	no	no	no	no	yes	no	yes	yes		<p>Why does this male have his arms in this position?</p> <table border="1"><tr><td>balance</td><td>angry</td></tr><tr><td>for balance</td><td>he's carrying bags</td></tr><tr><td>for balance</td><td>hug</td></tr></table>	balance	angry	for balance	he's carrying bags	for balance	hug		<p>How many people are wearing an orange shirt?</p> <table border="1"><tr><td>3</td><td>1</td></tr><tr><td>3</td><td>3</td></tr><tr><td>3</td><td>3</td></tr></table>	3	1	3	3	3	3		<p>Which player on the field head-butted the ball?</p> <table border="1"><tr><td>18</td><td>1</td><td>18</td></tr><tr><td>18</td><td>in front of goal</td><td>13</td></tr><tr><td>player on left</td><td>number 22</td><td>18</td></tr></table>	18	1	18	18	in front of goal	13	player on left	number 22	18
no	no	no																																			
no	no	yes																																			
no	yes	yes																																			
balance	angry																																				
for balance	he's carrying bags																																				
for balance	hug																																				
3	1																																				
3	3																																				
3	3																																				
18	1	18																																			
18	in front of goal	13																																			
player on left	number 22	18																																			
	<p>What are the people throwing?</p> <table border="1"><tr><td>frisbee</td><td>balls</td></tr><tr><td>frisbee</td><td>rice</td></tr><tr><td>frisbee</td><td>frisbee</td></tr></table>	frisbee	balls	frisbee	rice	frisbee	frisbee		<p>Are the clouds high in the sky?</p> <table border="1"><tr><td>yes</td><td>no</td></tr><tr><td>yes</td><td>no</td></tr><tr><td>yes</td><td>yes</td></tr></table>	yes	no	yes	no	yes	yes		<p>Is this a trained elephant?</p> <table border="1"><tr><td>yes</td><td>yes</td><td>yes</td></tr><tr><td>yes</td><td>yes</td><td>yes</td></tr><tr><td>yes</td><td>yes</td><td>yes</td></tr></table>	yes	yes	yes	yes	yes	yes	yes	yes	yes		<p>What number is on the girl in black?</p> <table border="1"><tr><td>18</td><td>1</td></tr><tr><td>18</td><td>4</td></tr><tr><td>18</td><td>8</td></tr></table>	18	1	18	4	18	8			
frisbee	balls																																				
frisbee	rice																																				
frisbee	frisbee																																				
yes	no																																				
yes	no																																				
yes	yes																																				
yes	yes	yes																																			
yes	yes	yes																																			
yes	yes	yes																																			
18	1																																				
18	4																																				
18	8																																				
	<p>What is the woman carrying?</p> <table border="1"><tr><td>umbrella</td><td>phone</td></tr><tr><td>umbrella</td><td>purse</td></tr><tr><td>umbrella</td><td>suitcase</td></tr></table>	umbrella	phone	umbrella	purse	umbrella	suitcase		<p>What is the type of meat under the bread?</p> <table border="1"><tr><td>beef</td><td>bologna</td></tr><tr><td>roast beef</td><td>ham</td></tr><tr><td>roast beef</td><td>tan</td></tr></table>	beef	bologna	roast beef	ham	roast beef	tan		<p>Is it sunny in this picture?</p> <table border="1"><tr><td>yes</td><td>yes</td><td>yes</td></tr><tr><td>yes</td><td>yes</td><td>yes</td></tr><tr><td>yes</td><td>yes</td><td>yes</td></tr></table>	yes	yes	yes	yes	yes	yes	yes	yes	yes		<p>How many stories is this home?</p> <table border="1"><tr><td>1</td><td>2</td></tr><tr><td>1</td><td>3</td></tr><tr><td>2</td><td>3</td></tr></table>	1	2	1	3	2	3			
umbrella	phone																																				
umbrella	purse																																				
umbrella	suitcase																																				
beef	bologna																																				
roast beef	ham																																				
roast beef	tan																																				
yes	yes	yes																																			
yes	yes	yes																																			
yes	yes	yes																																			
1	2																																				
1	3																																				
2	3																																				
	<p>What is the most colorful object in the picture?</p> <table border="1"><tr><td>umbrella</td><td>art</td></tr><tr><td>umbrella</td><td>flower</td></tr><tr><td>umbrella</td><td>flowers</td></tr></table>	umbrella	art	umbrella	flower	umbrella	flowers		<p>How much fluid is in the bottom of the bowl?</p> <table border="1"><tr><td>0</td><td>3 ounces</td></tr><tr><td>1 oz</td><td>a lot</td></tr><tr><td>little</td><td>some</td></tr></table>	0	3 ounces	1 oz	a lot	little	some		<p>How many bikes on the floor?</p> <table border="1"><tr><td>2</td><td>3</td></tr><tr><td>2</td><td>3 bikes</td></tr><tr><td>2</td><td>4</td></tr></table>	2	3	2	3 bikes	2	4		<p>Is the table large enough for 10 average people to eat at?</p> <table border="1"><tr><td>no</td><td>no</td></tr><tr><td>no</td><td>yes</td></tr><tr><td>no</td><td>yes</td></tr></table>	no	no	no	yes	no	yes						
umbrella	art																																				
umbrella	flower																																				
umbrella	flowers																																				
0	3 ounces																																				
1 oz	a lot																																				
little	some																																				
2	3																																				
2	3 bikes																																				
2	4																																				
no	no																																				
no	yes																																				
no	yes																																				
	<p>What is the guy doing as he sits on the bench?</p> <table border="1"><tr><td>phone</td><td>reading</td></tr><tr><td>taking picture</td><td>reading</td></tr><tr><td>taking picture with phone</td><td>smokes</td></tr></table>	phone	reading	taking picture	reading	taking picture with phone	smokes		<p>What color are his shoes?</p> <table border="1"><tr><td>blue</td><td>black</td></tr><tr><td>blue</td><td>black</td></tr><tr><td>blue</td><td>brown</td></tr></table>	blue	black	blue	black	blue	brown		<p>What is the horse missing to be able to ride it?</p> <table border="1"><tr><td>saddle</td><td>saddle</td></tr><tr><td>saddle</td><td>saddle</td></tr><tr><td>saddle</td><td>saddle</td></tr></table>	saddle	saddle	saddle	saddle	saddle	saddle		<p>What is the animal in the water?</p> <table border="1"><tr><td>dog</td><td>duck</td></tr><tr><td>dog</td><td>duck</td></tr><tr><td>dog</td><td>guppy</td></tr></table>	dog	duck	dog	duck	dog	guppy						
phone	reading																																				
taking picture	reading																																				
taking picture with phone	smokes																																				
blue	black																																				
blue	black																																				
blue	brown																																				
saddle	saddle																																				
saddle	saddle																																				
saddle	saddle																																				
dog	duck																																				
dog	duck																																				
dog	guppy																																				
	<p>What red objects in front are almost covered by snow?</p> <table border="1"><tr><td>meters</td><td>car</td></tr><tr><td>parking meters</td><td>cars</td></tr><tr><td>parking meters</td><td>shoes</td></tr></table>	meters	car	parking meters	cars	parking meters	shoes		<p>Is it winter?</p> <table border="1"><tr><td>yes</td><td>no</td></tr><tr><td>yes</td><td>yes</td></tr><tr><td>yes</td><td>yes</td></tr></table>	yes	no	yes	yes	yes	yes		<p>Does the car have a license plate?</p> <table border="1"><tr><td>yes</td><td>yes</td></tr><tr><td>yes</td><td>yes</td></tr><tr><td>yes</td><td>yes</td></tr></table>	yes	yes	yes	yes	yes	yes		<p>How many people are present?</p> <table border="1"><tr><td>15</td><td>2</td></tr><tr><td>15</td><td>3</td></tr><tr><td>15</td><td>3</td></tr></table>	15	2	15	3	15	3						
meters	car																																				
parking meters	cars																																				
parking meters	shoes																																				
yes	no																																				
yes	yes																																				
yes	yes																																				
yes	yes																																				
yes	yes																																				
yes	yes																																				
15	2																																				
15	3																																				
15	3																																				
	<p>Is this photo taken in Antarctica?</p> <table border="1"><tr><td>no</td><td>no</td></tr><tr><td>no</td><td>yes</td></tr><tr><td>no</td><td>yes</td></tr></table>	no	no	no	yes	no	yes		<p>Overcast or sunny?</p> <table border="1"><tr><td>overcast</td><td>overcast</td></tr><tr><td>overcast</td><td>overcast</td></tr><tr><td>overcast</td><td>sunny</td></tr></table>	overcast	overcast	overcast	overcast	overcast	sunny		<p>Could the truck have a camper?</p> <table border="1"><tr><td>yes</td><td>yes</td></tr><tr><td>yes</td><td>yes</td></tr><tr><td>yes</td><td>yes</td></tr></table>	yes	yes	yes	yes	yes	yes		<p>What is on the ground?</p> <table border="1"><tr><td>snow</td><td>dirt</td></tr><tr><td>snow</td><td>dirt</td></tr><tr><td>snow</td><td>mud</td></tr></table>	snow	dirt	snow	dirt	snow	mud						
no	no																																				
no	yes																																				
no	yes																																				
overcast	overcast																																				
overcast	overcast																																				
overcast	sunny																																				
yes	yes																																				
yes	yes																																				
yes	yes																																				
snow	dirt																																				
snow	dirt																																				
snow	mud																																				
	<p>Is the picture hanging straight?</p> <table border="1"><tr><td>no</td><td>no</td></tr><tr><td>yes</td><td>yes</td></tr><tr><td>yes</td><td>yes</td></tr></table>	no	no	yes	yes	yes	yes		<p>How many cabinets are on the piece of furniture?</p> <table border="1"><tr><td>4</td><td>3</td></tr><tr><td>4</td><td>3</td></tr><tr><td>4</td><td>6</td></tr></table>	4	3	4	3	4	6		<p>Why is the woman holding an umbrella?</p> <table border="1"><tr><td>sunny</td><td>it's raining</td></tr><tr><td>to block sun</td><td>it's raining</td></tr><tr><td>uncertain</td><td>to stay dry</td></tr></table>	sunny	it's raining	to block sun	it's raining	uncertain	to stay dry		<p>Does the man have a backpack?</p> <table border="1"><tr><td>yes</td><td>no</td></tr><tr><td>yes</td><td>yes</td></tr><tr><td>yes</td><td>yes</td></tr></table>	yes	no	yes	yes	yes	yes						
no	no																																				
yes	yes																																				
yes	yes																																				
4	3																																				
4	3																																				
4	6																																				
sunny	it's raining																																				
to block sun	it's raining																																				
uncertain	to stay dry																																				
yes	no																																				
yes	yes																																				
yes	yes																																				
	<p>What type of trees are here?</p> <table border="1"><tr><td>palm</td><td>ash</td></tr><tr><td>palm</td><td>oak</td></tr><tr><td>palm</td><td>pine</td></tr></table>	palm	ash	palm	oak	palm	pine		<p>Is the skateboard airborne?</p> <table border="1"><tr><td>yes</td><td>no</td></tr><tr><td>yes</td><td>yes</td></tr><tr><td>yes</td><td>yes</td></tr></table>	yes	no	yes	yes	yes	yes		<p>Is this person trying to hit a ball?</p> <table border="1"><tr><td>yes</td><td>yes</td><td>yes</td></tr><tr><td>yes</td><td>yes</td><td>yes</td></tr><tr><td>yes</td><td>yes</td><td>yes</td></tr></table>	yes	yes	yes	yes	yes	yes	yes	yes	yes		<p>What is the person hitting the ball with?</p> <table border="1"><tr><td>frisbee</td><td>bat</td></tr><tr><td>racket</td><td>bat</td></tr><tr><td>round paddle</td><td>racket</td></tr></table>	frisbee	bat	racket	bat	round paddle	racket			
palm	ash																																				
palm	oak																																				
palm	pine																																				
yes	no																																				
yes	yes																																				
yes	yes																																				
yes	yes	yes																																			
yes	yes	yes																																			
yes	yes	yes																																			
frisbee	bat																																				
racket	bat																																				
round paddle	racket																																				

Fig. 18: Random examples of questions (black), answers given when looking at the image (green), and answers given when not looking at the image (blue) for numerous representative examples of the real image dataset.

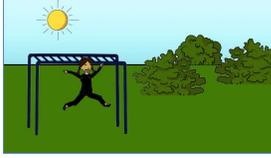
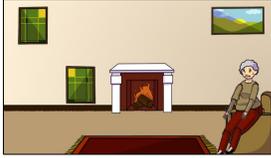
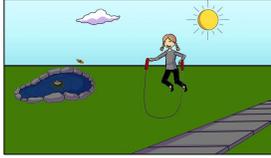
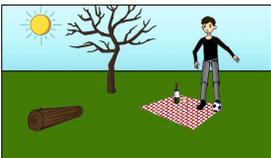
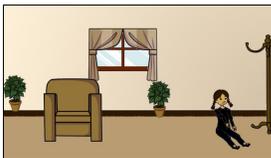
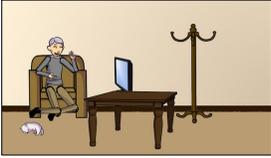
			
Who is holding the football? man boy man girl	What is the woman doing? sitting reading sitting watching tv	How many bushes are in the background? 3 3 3	What color is the bike? orange orange orange blue pink red
How is the weather? cool and sunny mostly sunny partly cloudy	Who is having tea? lady woman woman woman	What is the girl doing? playing playing playing	Is the man injured? no no no no yes
			
What is the dog looking at? ball soccerball cat tree	What color is the scooter? red red red red yellow	What part of the chair is the lady sitting on? arm arm arm arm seat seat	What is the little girl playing with? jump rope jump rope jump rope doll dolls teddy bear
Will the boy play with the dog? yes yes yes yes	How many turtles? 2 2 2 2 15	Is the woman sad? her cat died yes yes no no	What is in the pond? frog lily pad lily pad fish fish turtle
			
Are there leaves in the tree? no no no yes yes	What color is the book the woman is reading? blue blue blue green red	What is the girl sitting on? floor floor floor bench chair rock	What are the boy and girl sitting on? seesaw see saw teeter-totter bench couch
What is under the mans left foot? ball soccer ball tablecloth dollar grass ground	Is the lady reading? no yes yes no yes	What is the girl doing? sitting on floor sit ups dancing singing sleeping	What geometric shape is the base of the seesaw? triangle triangle triangle triangle triangle
			
Is the man happy? yes yes yes yes yes	How many cats are sleeping on the rug? 1 2 4 1 2	Is the sun shining? yes yes yes no yes yes	Does the man have a good heart? no yes yes yes yes
Is there an animal in the picture? yes yes yes yes	What color is the dog on the left? brown brown and white tan and white brown brown	What is in the pond? duck duck duck ducks fish fish	How many rabbits are there? 4 4 4 4 4
			
How many different kinds of fruits are available? 2 2 2 3 4 7	Is the cat chasing the mouse? yes yes yes yes yes	Is the man young or old? old old oldish old old	Is it a warm night? no no no no yes
Which objects needs 2 people in order to work? hands seesaw teeter-totter bandsaw, firehose jumprope seesaw	Is the man sad? no yes yes yes yes	Which grill is the man using? 1 on left left left barbeque gas left 1	Is the man happy? my best guess is happy yes yes no yes
			
How many windows are in this room? 2 2 2 2 4 8	Is the woman standing? no no no yes yes	What color is the plant on the left? green green green green red	How many books are in the shelf? 3 3 3 3 9 23
Is she waiting on someone? yes yes no no yes	What is beside the chair? dog dog plant cat table table	Why is the woman eating a salad rather than pizza? dieting on diet she likes salad dieting overweight she's losing weight	What is the person holding? book book notebook phone phone tablet

Fig. 19: Random examples of questions (black), answers given when looking at the image (green), and answers given when not looking at the image (blue) for numerous representative examples of the abstract scene dataset.

Q: What is the horse missing to be able to ride it?

(a) yes (b) no (c) 1 (d) 2 (e) 3 (f) 4
 (g) black (h) white (i) blue (j) red
 (k) nova (l) saddle (m) 1 week (n) yuy 551
 (o) tree (p) carriage (q) falling (r) cheese nuts spinach

Q: What shape is the building on the right?

(a) yes (b) no (c) 1 (d) 2 (e) 3 (f) 4
 (g) black (h) white (i) blue (j) red
 (k) steeple (l) soccerfield (m) rectangular (n) pyramid
 (o) cherry (p) square (q) bear and lizard (r) triangle

Q: Is the woman on the back of the bicycle pedaling?

(a) yes (b) no (c) 1 (d) 2 (e) 3 (f) 4
 (g) black (h) white (i) blue (j) red
 (k) relaxed (l) desk (m) in snow (n) marinara
 (o) tv show (p) tarach (q) shirt (r) bunny rabbit

Q: Why is the woman holding an umbrella?

(a) yes (b) no (c) 1 (d) 2 (e) 3 (f) 4
 (g) black (h) white (i) blue (j) red
 (k) it's raining (l) fountain (m) to block sun (n) rain
 (o) uncertain (p) to stay dry (q) orange juice expresso (r) sunny crema water

Q: Is the man happy?

(a) yes (b) no (c) 1 (d) 2 (e) 3 (f) 4
 (g) black (h) white (i) blue (j) red
 (k) unsure (l) butterflys (m) sitting (n) jeans
 (o) sad (p) happy (q) sunsetting (r) koi fish and lilypad over water

Q: Is there an animal in the picture?

(a) yes (b) no (c) 1 (d) 2 (e) 3 (f) 4
 (g) black (h) white (i) blue (j) red
 (k) 2 years (l) v (m) climbing onto sofa (n) 0
 (o) walked (p) scared of dog (q) is tree grown (r) mouth

Q: How many people are in the picture?

(a) yes (b) no (c) 1 (d) 2 (e) 3 (f) 4
 (g) black (h) white (i) blue (j) red
 (k) 5 (l) getting hit (m) 10 (n) sitting
 (o) kababs (p) trout (q) 77 (r) scratching his forehead

Q: What type of ball is hiding behind the tree?

(a) yes (b) no (c) 1 (d) 2 (e) 3 (f) 4
 (g) black (h) white (i) blue (j) red
 (k) sitting (l) earl gray (m) volleyball (n) soccer ball
 (o) acting (p) soccer (q) ball hiding behind tree is big ball (r) tennis

Q: What is on the ground?

(a) yes (b) no (c) 1 (d) 2 (e) 3 (f) 4
 (g) black (h) white (i) blue (j) red
 (k) dough (l) snow (m) toys (n) behind weirdo
 (o) dirt (p) shows (q) mud (r) microsoft

Q: Does the man have a backpack?

(a) yes (b) no (c) 1 (d) 2 (e) 3 (f) 4
 (g) black (h) white (i) blue (j) red
 (k) pizza cutter (l) ottoman has chrome finish (m) white black turquoise (n) silver
 (o) chrome (p) on surfboard in ocean (q) red vehicle w lights on beyond dark vehicle car lenth back (r) gravey

Q: Is the picture hanging straight?

(a) yes (b) no (c) 1 (d) 2 (e) 3 (f) 4
 (g) black (h) white (i) blue (j) red
 (k) straight (l) tooth brush (m) blue green (n) wrinkled
 (o) that it s on (p) highway (q) nope (r) turkey cheese

Q: How many cabinets are on the piece of furniture?

(a) yes (b) no (c) 1 (d) 2 (e) 3 (f) 4
 (g) black (h) white (i) blue (j) red
 (k) dumb (l) chest (m) herding competition (n) 6
 &dumber (o) 5 (p) 0 (q) 10-speed (r) ja8985

Q: Which objects needs 2 people in order for it to work?

(a) yes (b) no (c) 1 (d) 2 (e) 3 (f) 4
 (g) black (h) white (i) blue (j) red
 (k) plants (l) sitting (m) seesaw (n) her pants
 (o) hands (p) he s home (q) bandsaw firehose (r) 0

Q: What sport are they playing?

(a) yes (b) no (c) 1 (d) 2 (e) 3 (f) 4
 (g) black (h) white (i) blue (j) red
 (k) soccer (l) lacrosse (m) no he is not (n) sitting
 (o) frisbee (p) football (q) 0 (r) cycle

Q: What is the woman reaching for?

(a) yes (b) no (c) 1 (d) 2 (e) 3 (f) 4
 (g) black (h) white (i) blue (j) red
 (k) ties (l) glass (m) wine (n) sitting
 (o) question (p) dish (q) remote (r) door handle

Q: Who is playing with the dog?

(a) yes (b) no (c) 1 (d) 2 (e) 3 (f) 4
 (g) black (h) white (i) blue (j) red
 (k) 0 (l) boy (m) hopping bunny (n) ball
 (o) his owner (p) sitting (q) lilly pads (r) man

Fig. 20: Random examples of multiple-choice questions for numerous representative examples of the real and abstract scene dataset.

REFERENCES

- [1] S. Antol, C. L. Zitnick, and D. Parikh. Zero-Shot Learning via Visual Abstraction. In *ECCV*, 2014. [2](#), [3](#), [10](#)
- [2] J. P. Bigham, C. Jayant, H. Ji, G. Little, A. Miller, R. C. Miller, R. Miller, A. Tatarowicz, B. White, S. White, and T. Yeh. VizWiz: Nearly Real-time Answers to Visual Questions. In *User Interface Software and Technology*, 2010. [1](#), [2](#), [8](#)
- [3] K. Bollacker, C. Evans, P. Paritosh, T. Sturge, and J. Taylor. Freebase: A Collaboratively Created Graph Database for Structuring Human Knowledge. In *International Conference on Management of Data*, 2008. [2](#)
- [4] A. Carlson, J. Betteridge, B. Kisiel, B. Settles, E. R. H. Jr., and T. M. Mitchell. Toward an Architecture for Never-Ending Language Learning. In *AAAI*, 2010. [2](#)
- [5] X. Chen, H. Fang, T.-Y. Lin, R. Vedantam, S. Gupta, P. Dollár, and C. L. Zitnick. Microsoft COCO Captions: Data Collection and Evaluation Server. *arXiv preprint arXiv:1504.00325*, 2015. [3](#)
- [6] X. Chen, A. Shrivastava, and A. Gupta. NEIL: Extracting Visual Knowledge from Web Data. In *ICCV*, 2013. [2](#)
- [7] X. Chen and C. L. Zitnick. Mind’s Eye: A Recurrent Visual Representation for Image Caption Generation. In *CVPR*, 2015. [1](#), [2](#)
- [8] G. Coppersmith and E. Kelly. Dynamic wordclouds and vennclouds for exploratory data analysis. In *ACL Workshop on Interactive Language Learning and Visualization*, 2014. [9](#)
- [9] J. Deng, A. C. Berg, and L. Fei-Fei. Hierarchical Semantic Indexing for Large Scale Image Retrieval. In *CVPR*, 2011. [2](#)
- [10] J. Donahue, L. A. Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, and T. Darrell. Long-term Recurrent Convolutional Networks for Visual Recognition and Description. In *CVPR*, 2015. [1](#), [2](#)
- [11] D. Elliott and F. Keller. Comparing Automatic Evaluation Measures for Image Description. In *ACL*, 2014. [1](#)
- [12] A. Fader, L. Zettlemoyer, and O. Etzioni. Paraphrase-Driven Learning for Open Question Answering. In *ACL*, 2013. [2](#)
- [13] A. Fader, L. Zettlemoyer, and O. Etzioni. Open Question Answering over Curated and Extracted Knowledge Bases. In *International Conference on Knowledge Discovery and Data Mining*, 2014. [2](#)
- [14] H. Fang, S. Gupta, F. N. Iandola, R. Srivastava, L. Deng, P. Dollár, J. Gao, X. He, M. Mitchell, J. C. Platt, C. L. Zitnick, and G. Zweig. From Captions to Visual Concepts and Back. In *CVPR*, 2015. [1](#), [2](#)
- [15] A. Farhadi, M. Hejrati, A. Sadeghi, P. Young, C. Rashtchian, J. Hockenmaier, and D. Forsyth. Every Picture Tells a Story: Generating Sentences for Images. In *ECCV*, 2010. [2](#)
- [16] D. Geman, S. Geman, N. Hallonquist, and L. Younes. A Visual Turing Test for Computer Vision Systems. In *PNAS*, 2014. [1](#), [2](#)
- [17] J. Gordon and B. V. Durme. Reporting bias and knowledge extraction. In *Proceedings of the 3rd Workshop on Knowledge Extraction, at CIKM 2013*, 2013. [9](#)
- [18] S. Guadarrama, N. Krishnamoorthy, G. Malkarnenkar, S. Venugopalan, R. Mooney, T. Darrell, and K. Saenko. YouTube2Text: Recognizing and Describing Arbitrary Activities Using Semantic Hierarchies and Zero-Shot Recognition. In *ICCV*, December 2013. [2](#)
- [19] M. Hodosh, P. Young, and J. Hockenmaier. Framing Image Description as a Ranking Task: Data, Models and Evaluation Metrics. *JAIR*, 2013. [1](#)
- [20] A. Karpathy and L. Fei-Fei. Deep Visual-Semantic Alignments for Generating Image Descriptions. In *CVPR*, 2015. [1](#), [2](#)
- [21] S. Kazemzadeh, V. Ordonez, M. Matten, and T. L. Berg. ReferItGame: Referring to Objects in Photographs of Natural Scenes. In *EMNLP*, 2014. [2](#)
- [22] R. Kiros, R. Salakhutdinov, and R. S. Zemel. Unifying Visual-Semantic Embeddings with Multimodal Neural Language Models. *TACL*, 2015. [1](#), [2](#)
- [23] C. Kong, D. Lin, M. Bansal, R. Urtasun, and S. Fidler. What Are You Talking About? Text-to-Image Coreference. In *CVPR*, 2014. [2](#)
- [24] A. Krizhevsky, I. Sutskever, and G. E. Hinton. ImageNet Classification with Deep Convolutional Neural Networks. In *NIPS*, 2012. [2](#), [7](#)
- [25] G. Kulkarni, V. Premraj, S. L. Sagnik Dhar and, Y. Choi, A. C. Berg, and T. L. Berg. Baby Talk: Understanding and Generating Simple Image Descriptions. In *CVPR*, 2011. [2](#)
- [26] D. B. Lenat and R. V. Guha. *Building Large Knowledge-Based Systems; Representation and Inference in the Cyc Project*. Addison-Wesley Longman Publishing Co., Inc., 1989. [2](#)
- [27] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft COCO: Common Objects in Context. In *ECCV*, 2014. [2](#), [3](#), [4](#), [12](#)
- [28] X. Lin and D. Parikh. Don’t Just Listen, Use Your Imagination: Leveraging Visual Common Sense for Non-Visual Tasks. In *CVPR*, 2015. [1](#)
- [29] H. Liu and P. Singh. ConceptNet — A Practical Commonsense Reasoning Tool-Kit. *BT Technology Journal*, 2004. [2](#)
- [30] M. Malinowski and M. Fritz. A Multi-World Approach to Question Answering about Real-World Scenes based on Uncertain Input. In *NIPS*, 2014. [1](#), [2](#)
- [31] J. Mao, W. Xu, Y. Yang, J. Wang, and A. L. Yuille. Explain Images with Multimodal Recurrent Neural Networks. *CoRR*, abs/1410.1090, 2014. [1](#), [2](#)
- [32] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean. Distributed Representations of Words and Phrases and their Compositionality. In *NIPS*, 2013. [4](#)
- [33] M. Mitchell, J. Dodge, A. Goyal, K. Yamaguchi, K. Stratos, X. Han, A. Mensch, A. Berg, T. L. Berg, and H. Daume III. Midge: Generating Image Descriptions From Computer Vision Detections. In *ACL*, 2012. [2](#)
- [34] M. Mitchell, K. van Deemter, and E. Reiter. Attributes in visual reference. In *PRE-CogSci*, 2013. [9](#)
- [35] M. Mitchell, K. Van Deemter, and E. Reiter. Generating Expressions that Refer to Visible Objects. In *HLT-NAACL*, 2013. [2](#)
- [36] V. Ramanathan, A. Joulin, P. Liang, and L. Fei-Fei. Linking People with “Their” Names using Coreference Resolution. In *ECCV*, 2014. [2](#)
- [37] M. Richardson, C. J. Burges, and E. Renshaw. MCTest: A Challenge Dataset for the Open-Domain Machine Comprehension of Text. In *EMNLP*, 2013. [1](#), [2](#)
- [38] M. Rohrbach, W. Qiu, I. Titov, S. Thater, M. Pinkal, and B. Schiele. Translating Video Content to Natural Language Descriptions. In *ICCV*, 2013. [2](#)
- [39] K. Toutanova, D. Klein, C. D. Manning, and Y. Singer. Feature-rich part-of-speech tagging with a cyclic dependency network. In *ACL*, 2003. [8](#)
- [40] K. Tu, M. Meng, M. W. Lee, T. E. Choe, and S. C. Zhu. Joint Video and Text Parsing for Understanding Events and Answering Queries. *IEEE MultiMedia*, 2014. [1](#), [2](#)
- [41] R. Vedantam, C. L. Zitnick, and D. Parikh. CIDER: Consensus-based Image Description Evaluation. In *CVPR*, 2015. [1](#)
- [42] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan. Show and Tell: A Neural Image Caption Generator. In *CVPR*, 2015. [1](#), [2](#)
- [43] J. Weston, A. Bordes, S. Chopra, and T. Mikolov. Towards AI-Complete Question Answering: A Set of Prerequisite Toy Tasks. *CoRR*, abs/1502.05698, 2015. [2](#)
- [44] C. L. Zitnick and D. Parikh. Bringing Semantics Into Focus Using Visual Abstraction. In *CVPR*, 2013. [2](#), [3](#)
- [45] C. L. Zitnick, D. Parikh, and L. Vanderwende. Learning the Visual Interpretation of Sentences. In *ICCV*, 2013. [2](#), [3](#)
- [46] C. L. Zitnick, R. Vedantam, and D. Parikh. Adopting Abstract Images for Semantic Scene Understanding. *PAMI*, 2015. [2](#), [3](#)