

GQA: a new dataset for compositional question answering over real-world images

visualreasoning.net

Drew A. Hudson
Stanford University
450 Serra Mall, Stanford, CA 94305
dorarad@cs.stanford.edu

Christopher D. Manning
Stanford University
450 Serra Mall, Stanford, CA 94305
manning@cs.stanford.edu

Abstract

We introduce *GQA*, a new dataset for real-world visual reasoning and compositional question answering, seeking to address key shortcomings of previous VQA datasets. We have developed a strong and robust question engine that leverages scene graph structures to create 22M diverse reasoning questions, all come with functional programs that represent their semantics. We use the programs to gain tight control over the answer distribution and present a new tunable smoothing technique to mitigate language biases. Accompanying the dataset is a suite of new metrics that evaluate essential qualities such as consistency, grounding and plausibility. An extensive analysis is performed for baselines as well as state-of-the-art models, providing fine-grained results for different question types and topologies. Whereas a blind LSTM obtains mere 42.1%, and strong VQA models achieve 54.1%, human performance tops at 89.3%, offering ample opportunity for new research to explore. We strongly hope *GQA* will provide an enabling resource for the next generation of models with enhanced robustness, improved consistency, and deeper semantic understanding for images and language.

1. Introduction

It takes more than a smart guess to answer a good question. The ability to truly assimilate knowledge and use it to draw inferences is among the holy grails of artificial intelligence. A tangible form of this goal is embodied in the task of Visual Question Answering (VQA), where a system has to answer free-form questions by reasoning about presented images. The task demands a rich set of abilities as varied as object recognition, commonsense understanding and relation extraction, spanning both the visual and linguistic ends. In recent years, it has sparked a substantial interest throughout the research community, becoming ex-



Figure 1: Examples from the new GQA dataset for visual reasoning and compositional question answering:
Is the **white** to the right of the **green** **apple**?
What type of **fruit** in the image is **round**?
What color is the **fruit** on the right side, **red** or **green**?
Is there any **milk** in the **bowl** to the left of the **apple**?

tremely popular across the board, with a host of datasets being constructed [3, 10, 15, 43, 20] and numerous models being proposed [4, 40, 5, 9, 11].

The multi-modal nature of the task and the diversity of skills required to address different questions make VQA particularly challenging. Yet, designing a good test that will reflect its full qualities and complications may not be as trivial. In spite of the great strides the field recently made, it has been established through a series of studies that existing benchmarks suffer from critical vulnerabilities that render them highly unreliable in measuring the actual degree of visual understanding capacities [41, 10, 1, 7, 2, 13, 18].

Most notable flaw of existing benchmarks is the strong language priors displayed throughout the data [41, 10, 2] – indeed, most tomatoes are red and most tables are wooden. These in turn are exploited by VQA models, which become

heavily reliant upon statistical biases and tendencies, rather than on true scene understanding skills [1, 10, 15]. They memorize the precise answer distributions of different questions to handle them with relative ease, while only glancing over the provided images and at times not even considering them, let alone understanding their content [1, 7]. Consequentially, early benchmarks have led to an inflated sense of the state of visual scene understanding, severely diminishing their credibility [38].

Apart from the prevalent biases within the questions, current real-image VQA datasets suffer from multiple other issues and deficiencies: For one thing, they commonly use basic, non-compositional language which rarely require far beyond object recognition [34]. Second, the immense variability in potential ways to refer or describe objects and scenes make it particularly hard for systems to distill/capture and learn clear/unambiguous grounded semantics, a crucial element of cogent scene understanding. Finally, the lack of annotations regarding questions structure, type and content leave it difficult to identify and fix the root causes behind mistakes models make [15].

To address these shortcomings, while retaining the visual and semantic richness of real-world images, we introduce GQA, a new dataset for visual reasoning and compositional question answering. We have developed and carefully refined a robust *question engine*, which leverages **content**: information about objects, attributes and relations provided through the Visual Genome *Scene Graphs* [20], along with **structure**: a newly-created extensive linguistic grammar which couples hundreds of structural patterns and detailed lexical semantic resources, partly derived from the VQA dataset. Together, we combine them to generate over 22 million novel and diverse questions, all come with structured representations in the form of functional programs that specify their contents and semantics, and are visually grounded in the image scene graphs.

Many of the GQA questions involve varied reasoning skills, and multi-step inference in particular, standing in sharp contrast with existing real-image VQA datasets [3, 10, 43] which tend to have fairly simple questions from both linguistic and semantic perspectives [34]. We further use the associated functional representations to greatly reduce biases within the dataset and control for its question type composition, downsampling it to create a 1.7M-questions balanced dataset. Contrary to VQA2.0 [10], here we balance not only binary questions, but also open ones, by applying a tunable smoothing technique that makes the answers distribution for each question group more uniform, thereby enabling tight control over the dataset composition. Just like a well designed exam, our benchmark makes the educated guesses strategy far less rewarding, and demands instead more refined comprehension of both the visual and linguistic contents. At the same time, we recognize the im-

portance of research on robustness against biases, and so will provide both the dataset’s balanced and original unbalanced versions.

Along with the data, we have designed a suite of new metrics, which include consistency, validity, plausibility, grounding and distribution scores, to complement the standard accuracy measure commonly used in assessing method’s performance. Indeed, studies have shown that the accuracy metric alone does not account for a range of anomalous behaviors that models demonstrate, such as ignoring key question words or attending to irrelevant image regions [1, 7]. Other works have argued for the need to devise new evaluation measures and techniques to shed more light on systems’ inner workings [18, 35, 36, 17]. In fact, beyond providing new metrics, GQA can even directly support the development of more interpretable models, as it provides a sentence-long explanation that corroborates each answer, and further associates each word from both the questions and the responses with a visual pointer to the relevant region/s in the image, similar in nature to datasets by Yuke *et al.* [43], Park *et al.* [30], and Li *et al.* [22]. These in turn can serve as a strong supervision signal to train models with enhanced transparency and accessibility.

In the following, we delineate the design of our question engine, explain the multi-step question generation process, analyze the resultant dataset and compare it with existing benchmarks. We present and discuss the new metrics and use them to evaluate an array of baselines and state-of-the-art models, comparing them to human subjects and revealing a large gap in performance across multiple axes. Finally, we discuss models’ strengths and weakness discovered through the analysis on GQA, and propose potential research directions to overcome them.

GQA combines the best of both worlds, having clearly defined and crisp semantic representations on the one hand but enjoying the semantic and visual richness of real-world images on the other. Our three main contributions are (1) the GQA dataset as a resource for studying visual reasoning; (2) development of an effective method for generating a large number of semantically varied questions, which marries scene graph representations with computational linguistic methods; (3) new metrics for GQA, that allow for better assessment of system success and failure modes, as demonstrated through a comprehensive performance analysis of existing models on this task. We hope that the GQA dataset will provide fertile ground for the development of novel methods that push the boundaries of questions answering and visual reasoning.

2. Related Work

The last few years have witnessed tremendous progress in visual understanding in general and VQA in particular, as we move beyond classic perceptual tasks towards problems

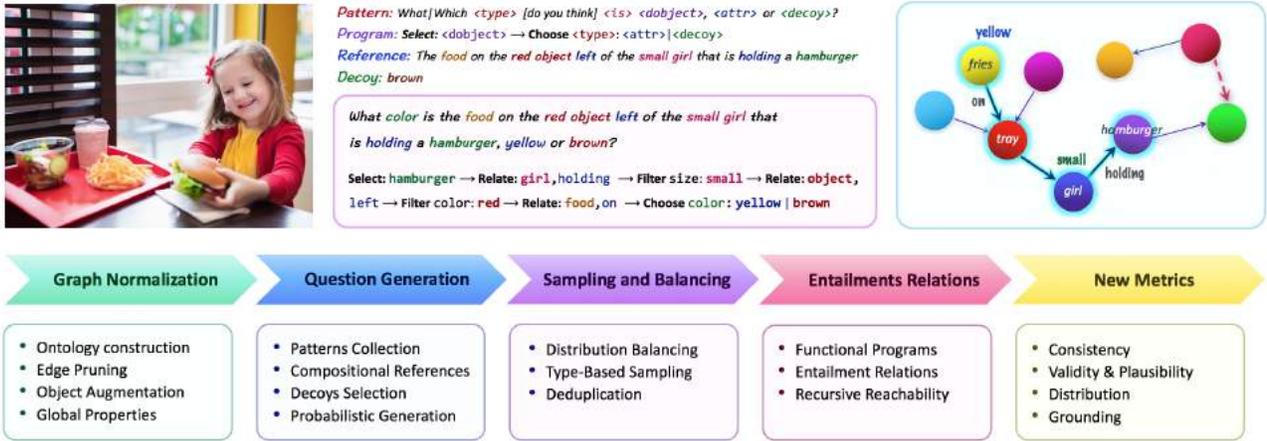


Figure 2: Overview of the GQA construction process. Given an image annotated with a Scene Graph of its objects, attributes and relations, we produce compositional questions by traversing the graph. Each question has both a standard natural-language form and a functional program representing its semantics. Please refer section 3.1 to for further detail.

that ask for high-level semantic understanding and integration of multiple modalities. However, as discussed in section 1, many of these benchmarks suffer from systematic biases, allowing models to circumvent the need for thorough visual understanding, and instead make use of the prevalent real-world language priors to predict plenty of answers with confidence. On the common VQA1.0 dataset, blind models achieve over 50% in accuracy without even considering the images whatsoever [3].

Initial attempts have been made to remedy this situation [10, 41, 2, 15], but they fall short in providing an adequate solution: Some approaches operate over constrained and synthetic images [41, 15], neglecting the realism and diversity natural photos provide. Meanwhile, Goyal *et al.* [10] associate most of the questions in VQA1.0 with a pair of similar pictures that result in different answers. While offering partial relief, this technique fails to address open questions, leaving their answer distribution largely unbalanced. In fact, since the method does not cover 29% of the questions, even within the binary ones biases still remain.¹

At the other extreme, Agrawal *et al.* [2] partition the questions into training and validation sets such that their respective answer distributions become intentionally dissimilar. While undoubtedly challenging, these adversarial settings penalize models, maybe unjustly, for learning salient properties of the training data. In the absence of other information, making an educated guess is actually the right choice – a valid and beneficial strategy pursued by machines

¹According to Goyal *et al.* [10], 22% of the original questions are left unpaired, and 9% of the paired ones get the same answer due to annotation errors. Indeed, baseline experiments reveal that 67% and 27% of the binary and open questions respectively are answered correctly by a blind model with no access to the input images.

and people alike [28, 6, 27]. While the ability to generalize in the face of change is certainly important, it is ancillary to the task of visual understanding in its purest form. Instead, what we essentially need is a fair but balanced test that is more resilient to such gaming strategies, as we strive to achieve with GQA.

In creating our dataset, we drew inspiration from the CLEVR task [15], which consists of compositional questions over synthetic images. However, its artificial nature and low diversity, with only 3 classes of objects and 12 different properties, makes it particularly vulnerable to memorization. In other words, its space is small enough that a model can easily learn an independent representation for each of the 96 combinations such as “large red sphere”, reducing its effective degree of compositionality. Conversely, GQA operates over real images and a large semantic space, making it much more realistic and challenging. Even though our questions are not natural as in other VQA datasets [10, 43], they display a broad vocabulary and diverse grammatical structures. They may serve in fact as a cleaner benchmark to assess models in a more controlled and comprehensive fashion, as discussed below. Specifically, our dataset builds on top of the *scene graph* annotations of Visual Genome [20], which is a crowdsourced dataset specifying the objects, attributes and relations present in 108K different images, all through natural, unconstrained language. Compared to synthetic datasets such as CLEVR, constructing a generation pipeline to cover such linguistic diversity and rich vocabulary entails unique challenges for GQA, as we further discuss in section 3.

Somewhat related to our work are several datasets created through templates [25, 17, 24], most of them for the purpose of data augmentation. However, they are either

small in scale [25] or use only a restricted set of objects and a handful of non-compositional templates [17, 24].² Neural alternatives to visual question generation have been recently proposed [29, 14, 42], but they aim at a quite different goal of creating “interesting and engaging” questions about the wider context of the image, *e.g.* subjective evoked feelings or speculative events that may lead or result from the depicted scenes [29]. On top of that, the neurally generated questions may actually be incorrect, irrelevant, or nonsensical. In contrast, here we are focusing on pertinent, factual and accurate questions with objective answers, seeking to create a challenging benchmark for the task of VQA.

3. GQA Dataset Creation

GQA is a new dataset for visual reasoning and compositional question answering over real-world images, designed to foster the development of models capable of advanced reasoning skills and improved scene understanding capabilities. By creating a balanced set of challenging questions over images, along with detailed annotations of the question and image semantics and a suite of new metrics, we allow comprehensive diagnosis of methods’ performance, and open the door for novel models with more transparent and coherent knowledge representation and reasoning.

Figure 2 provides a brief overview of the GQA components and generation process, and figure 3 offers multiple instances from the dataset. More examples are provided in figure 10. The dataset along with further information are available at visualreasoning.net.

3.1. Overview

The GQA dataset consists of 113K images and 22M questions of assorted types and varying compositionality degrees, measuring performance on an array of reasoning skills such as object and attribute recognition, transitive relation tracking, spatial reasoning, logical inference and comparisons.

The images, questions and corresponding answers are all accompanied by matching semantic representations: Each *image* is annotated with a dense **Scene Graph** [16, 20], representing the objects, attributes and relations it contains. Each *question* is associated with a **functional program** which lists the series of reasoning steps that have to be performed to arrive at the answer. Each *answer* is augmented with both textual and visual justifications, pointing to the relevant region within the image.

The structured representations offer multiple advantages, as they enable tight control over the question and answer distribution, discussed further in section 3.5, and facilitate

²According to [17, 24], 74% and 86% of their questions respectively are of the form “Is there X in the picture?”. In the latter, the answer to them is invariably “Yes”.

assessment of models’ performance along various axes, including question type, topology and semantic length (section 4.2 and section 10). In addition, they may support the development of more interpretable grounded models by serving as a strong supervision signal during training. Finally, they enable the design of a new consistency metric, as we can see in section 4.4.

We proceed by describing the four-step pipeline of the dataset construction: First, we thoroughly clean, normalize, consolidate and augment the Visual Genome scene graphs [20] linked to each image. Then, we traverse the graphs to collect information about objects and relations, which is then coupled with grammatical patterns gleaned from VQA2.0 [10] and sundry probabilistic grammar rules to produce semantically-rich and diverse set of questions. At the third stage, we use the underlying semantic forms to reduce biases in the conditional answer distribution, resulting in a balanced dataset that is more robust against shortcuts and guesses. Finally, we provide further detail about the questions’ functional representations, and explain how we utilize these to compute entailment between questions, which will be further used in section 4.4.

3.2. Scene Graph Normalization

Our starting point in creating the GQA dataset is the Visual Genome Scene Graph annotations³ [20] that cover 113k images from COCO [23] and Flickr [37].⁴ The scene graph serves as a formalized representation of the image: each node denotes an **object**, a visual entity within the image, like a person, an apple, grass or clouds. It is linked to a bounding box specifying its position and size, and is marked up with about 1-3 **attributes**, properties of the object: *e.g.* its color, shape, material or activity. The objects are connected by **relation** edges, representing actions (*i.e.* verbs), spatial relations (*e.g.* prepositions), and comparatives.

The scene graphs are annotated with unconstrained natural language. Our first goal is thus to convert the annotations into a clear and unambiguous semantic ontology.⁵ We begin by cleaning up the graphs vocabulary, removing stop words, fixing typos, consolidating synonyms and filtering rare or amorphous concepts.⁶ We then classify the vocabulary into predefined categories (*e.g.* *animals* and *fruits* for objects; *colors* and *materials* for attributes), using word embedding distances to get preliminary annotations, which are then followed by manual curation. This results in a class hierarchy

³We use the cleaner 1.4 version of the dataset, following Xu *et al.* [39]

⁴We expand the original Visual Genome dataset with 5k new scene graphs collected through crowdsourcing.

⁵Note that we cannot effectively use the wordnet annotations [26] used throughout the Visual Genome dataset since they are highly inaccurate, relating objects to irrelevant senses, *e.g.* accountant for a game controller, Cadmium for a CD, etc.

⁶During this stage we also address additional linguistic subtleties such as the use of noun phrases (“pocket watch”) and opaque compounds (“soft drink”, “hard disk”).



1. Is the **tray** on top of the **table** black or light brown? light brown
2. Are the **napkin** and the **cup** the same color? yes
3. Is the small **table** both oval and wooden? yes
4. Is the **syrup** to the left of the **napkin**? yes
5. Is there any **fruit** to the left of the **tray** the **cup** is on top of? yes
6. Are there any **cups** to the left of the **tray** that is on top of the **table**? no
7. Could this **room** be a living room? yes



1. Is there a **door** or a **window** that is open? no
2. Do you see any white **numbers** or **letters**? yes
3. What is the large **container** made of? cardboard
4. What **animal** is in the **box**? **bear**
5. Is there a **bag** right of the **bear**? no
6. Is there a **box** inside the plastic **bag**? no
7. What is the green **thing** on the right? **door**
8. What color is the **bear**? brown



1. Which side of the image is the **plate** on? right
2. Are there any **lamps** on the **desk** to the right of the **rug**? yes
3. What type of **furniture** are the **flowers** on, a **bed** or a **table**? **table**
4. Are there any **clocks** or **mirrors**? no
5. Are there any **chairs** to the right of the **lamp** on the **table**? yes
6. What is the dark **piece of furniture** to the right of the **rug** called? **cabinet**

Figure 3: Examples of questions from the GQA dataset.

over the scene graphs vocabulary, that we further augment with various semantic and linguistic features such as part of speech, voice, plurality and synonyms – information that will be used to create grammatically correct questions in further steps. Our final ontology contains 1740 objects, 620 attributes and 330 relations, grouped into a hierarchy that consists of 60 different categories and subcategories. Visualization of the ontology can be found in figure 9.

At the next step, we prune graph edges that sound unnatural or are otherwise inadequate to be incorporated within the questions to be generated, such as (*woman, in, shirt*), (*tail, attached to, giraffe*), or (*hand, hugging, bear*). We filter these triplets using a combination of category-based rules, n-gram frequencies [12], dataset co-occurrence statistics, and manual curation.

In order to generate correct and unambiguous questions, some cases will require us to validate the uniqueness or absence of an object. Visual Genome, while meant to be as exhaustive as possible, cannot guarantee full coverage (as it may be practically infeasible). Hence, in those cases we use object detectors [32], trained on visual genome with a low detection threshold, to conservatively confirm the object absence or uniqueness.⁷

Next, we augment the graph objects with absolute and relative positional information: objects appearing within margins, horizontal or vertical, are annotated accordingly. Object pairs for which we can safely determine horizontal

positional relations (*e.g.* one is to the left of the other), are annotated as well.⁸ We also annotate object pairs if they share the same color, material or shape. Finally, we enrich the graph with global information about the image location or weather, if these can be directly inferred from the objects it contains.

By the end of this stage, the resulting scene graphs have clean, unified, rich and unambiguous semantics for both the nodes and the edges.

3.3. The Question Engine

At the heart of our pipeline is the question engine, responsible for producing diverse, relevant and grammatical questions in varying degrees of compositionality. The generation process harnesses two resources: one is the scene graphs which fuel the engine with rich content – information about objects, attributes and relationships; the other is the structural patterns, a mold that shapes the content, casting it into a question.

Our engine operates over 524 patterns, spanning 117 question groups. Each group is associated with three components: (1) a functional program that represents its semantics; (2) A set of textual rephrases which express it in natural language, *e.g.* “What|Which <type> [do you think] <is> <theObject>?”; (3) A pair of short and long answers: *e.g.* <attribute> and “The <object> <is>

⁷An object uniqueness can be validated by confidently denying the existence of other same-class objects within the image.

⁸We do not annotate objects with vertical positional relations since these cannot be confidently determined from their bounding boxes – we cannot distinguish between cases of below/above and behind/in-front.

<attribute>.” respectively.⁹

We begin from a seed set of 250 manually constructed patterns, and extend it with 274 natural patterns derived from VQA1.0 [3] through anonymization of words from our ontology.¹⁰ To increase the question diversity, apart from using synonyms for objects and attributes, we incorporate probabilistic sections into the patterns, such as optional phrases $[x]$ and alternate expressions $(x|y)$, which get instantiated at random.

It is important to note that the patterns do not strictly limit the structure or depth of each question, but only outline their high-level form, as many of the template fields can be populated with nested compositional references. For instance, in the pattern above, we may replace $\langle \text{theObject} \rangle$ with ‘the apple to the left of the white refrigerator.

To achieve that compositionality, we compute for each object a set of candidate references, which can either be **direct**, e.g. *the bear*, *this animal*, or **indirect**, using modifiers, e.g. *the white bear*, *the bear on the left*, *the animal behind the tree*, *the bear that is wearing a coat*. Direct references are used when the uniqueness of the object can be confidently confirmed by object detectors, making the corresponding references unambiguous. Alternatively, we use indirect references, leading to multi-step questions as varied as *Who is looking at the animal that is wearing the red coat in front of the window?*, and thus greatly increasing the patterns’ effective flexibility. This is the key ingredient behind the automatic generation of compositional questions.

Finally, we compute a set of decoys for the scene graph elements. Indeed, some questions, such as negative ones or those that involve logical inference, pertain to the absence of an object or to an incorrect attribute. Examples include e.g. *Is the apple green?* for a red apple, or *Is the girl eating ice cream?* when she is in fact eating a cake. Given a triplet (s, r, o) , (e.g. *(girl, eating, cake)*) we select a distractor \hat{o} considering its likelihood to be in relation with s and its plausibility to co-occur in the context of the other objects in the depicted scene. Similar technique is applied in selecting attribute decoys (e.g. a green apple). While choosing distractors, we exclude from consideration candidates that we deem too similar (e.g. *pink* and *orange*), based on a manually defined list for each concept in the ontology.

Having all resources prepared: (1) the clean scene graphs, (2) the structural patterns, (3) the object references and (4) the decoys, we can proceed to generating the questions! We traverse the graph, and for each object, object-

⁹Note that the long answers can serve as textual justifications, especially for questions that require increased reasoning such as logical inference, where a question like “Is there a red apple in the picture?” may have the answer: “No, there is an apple, but it is green”

¹⁰For instance, a question-answer pair in VQA1.0 such as “What color is the apple? red” turns after anonymization into “What $\langle \text{type} \rangle$ $\langle \text{is} \rangle$ the $\langle \text{object} \rangle$? $\langle \text{attribute} \rangle$ ”.

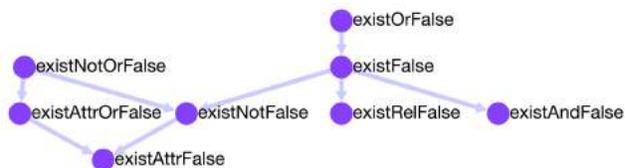


Figure 4: Examples of Entailment relations between different question types. Refer to section 3.4 for further detail.

attribute pair or subject-relation-object triplet, we produce relevant questions by instantiating a randomly selected question pattern, e.g. *What $\langle \text{type} \rangle$ is $\langle \text{theObject} \rangle$, $\langle \text{attribute} \rangle$ or $\langle \text{cAttribute} \rangle$?*, populating all the fields with the matching information, yielding, for example, the question: *What (color) (is) the (apple on the table), (red) or (green)?*. When choosing object references, we avoid selecting those that disclose the answer or repeat information, e.g. *What color is the red apple?* or *Which dessert sits besides the apple to the left of the cake?*. We also avoid asking about relations that tend to have multiple instances for the same object, e.g. asking what object is on the table, as they may be multiple valid answers.

Finally, we use the linguistic features associated with each vocabulary item along with n-gram frequencies [12] to correctly resolve final grammatical subtleties and further increase the questions linguistic variance, ironing out the last kinks and twists of the question engine.¹¹

By the end of this stage, we obtain a diverse set of 22M interesting, challenging and grammatical questions, pertaining to each and every aspect of the image.

3.4. Functional Representation and Entailment

Each question pattern is associated with a structured representation in the form of a functional program. For instance, the question *What color is the apple on the white table?* is semantically equivalent to the following program: *select: table \rightarrow filter: white \rightarrow relate(subject,on): apple \rightarrow query: color*. As we can see, these programs are composed of atomic operations such as object selection, traversal along a relation edge or an attribute verification, which are in turn concatenated together to create challenging reasoning questions.

The semantic unambiguous representations offer multiple advantages over free form unrestricted questions. For one thing, they enable comprehensive assessment of methods by dissecting their performance along different axes of question textual and semantic lengths, type and topology, thus facilitating the diagnosis of their success and failure

¹¹The considered nuances include determining prepositions, choosing articles and selecting the person for verbs and pronouns. Among other adjustments performed randomly, are changes to verb tense, use of contractions or apostrophes, and minor rearrangements in prepositions order.

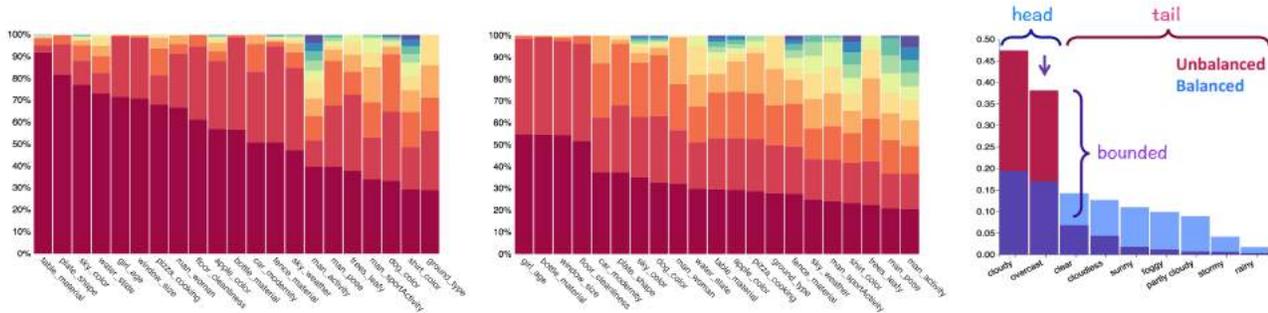


Figure 5: Visualization of the balancing process. **Left:** The conditional answer distribution before the balancing. We show the top 10 answers for a selection of question groups, where the column height corresponds to the relative frequency of each answer. We can see that the distributions are heavily biased. **Middle:** The distribution after balancing, more uniform and with heavier tails, while intentionally retaining the original real-world tendencies up to a tunable degree. **Right:** An illustration of the balancing process. Please refer to section 8 for further details and visualizations.

modes (section 4.2 and section 10). Second, they aid us in balancing the dataset distribution, mitigating its language priors and guarding against educated guesses (section 3.5). Finally, they allow us to identify entailment and equivalence relations between different questions: knowing the answer to the question *What color is the apple?* allows a coherent learner to infer the answer to the questions *Is the apple red?* *Is it green?* etc. The same goes especially for questions that involve logical inference like *or* and *and* operations or spatial reasoning, e.g. *left* and *right*. Please refer to figure 4 and figure 15 for entailment examples.

As further discussed in section 4.4, this entailment property can be used to measure the coherence and consistency of the models, shedding new light on their inner workings, compared to the widespread but potentially misleading accuracy metric. We define direct entailment relations between the various functional programs and use these to recursively compute all the questions that can be entailed from a given source. A complete catalog of the functions, their associated question types, and the entailment relations between them is provided in table 3 and figure 15.

3.5. Sampling and Balancing

One of the main issues of existing VQA datasets is the prevalent question-conditional biases that allow learners to make educated guesses without truly understanding the presented images, as explained in section 1. However, precise representation of the questions’ semantics can allow tighter control over these biases, having potential to greatly alleviate the problem. We leverage this observation and use the functional programs attached to each question to smooth out the answer distribution.

Given a question’s semantic program, we derive two labels, global and local: The global label assigns the question to its answer type, e.g. *color* for *What color is the apple?*.

The local label further considers the main subject/s of the question, e.g. *apple-color* or *table-material*. We use these labels to partition the questions into groups, and smooth the answer distribution of each group within the two levels of granularity, first globally, and then locally.

For each group, we first compute its answer distribution P , and then sort the answers based on their frequency within the group. Then, we downsample the questions (formally, using rejection-sampling) to fit a smoother answer distribution Q derived through the following procedure: We iterate over the answers in decreasing frequency order, and reweight P ’s head up to the current iteration to make it more comparable to the tail size. While repeating this operation as we go through the answers, iteratively “moving” probability from the head into the tail [33], we also maintain minimum and maximum ratios between each pair of subsequent answers (sorted by frequency). This ensures that the relative frequency-based answer ranking stays the same.

The main advantage of this scheme is that it retains the general real-world tendencies, smoothing them out up to a tunable degree to make the benchmark more challenging and less biased. Refer to figure 5 for a visualization and to section 8 for a precise depiction of the procedure. Since we perform this balancing in two granularity levels, the obtained answer distributions are made more uniform both locally and globally. Quantitatively, the entropy of the answer distribution is increased by 72%, confirming the success of this stage.

Finally, we downsample the questions based on their type to control the dataset type composition, and filter out redundant questions that are too semantically similar to existing ones. We split the dataset into 70% train, 10% validation, 10% test and 19% challenge, making sure that all the questions about a given image appear in the same split.

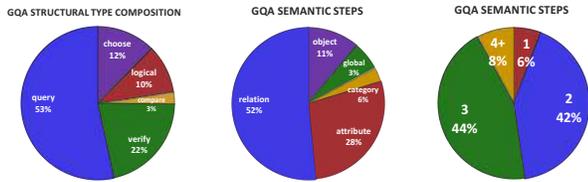


Figure 6: Dataset statistics: partitioned into structural types which indicate the final operation to be performed, semantic types, specifying the main subject of the question, and semantic length, the number of reasoning steps.

4. Analysis and Baseline experiments

In the following, we provide an analysis of the GQA dataset, perform a head-to-head comparison with the common VQA2.0 dataset [10], and evaluate the performance of baselines as well as state-of-the-art models. We introduce the new metrics that complement our dataset, provide quantitative results and discuss their implications and merits. To establish the diversity and realism of GQA questions, we then show test transfer performance between the GQA and VQA datasets. Finally, In section 10, we proceed with further diagnosis of the current top-performing model, MAC [11], evaluating it along multiple axes such as training-set size, question length and compositionality degree.

4.1. Dataset Analysis and Comparison

The GQA dataset includes 22,669,678 questions over 113,018 images. As figure 7 shows, the questions of varied lengths, longer than those of the VQA2.0 benchmark, alluding to their higher compositionality degree (figure 6). It has a vocabulary size of 3097 words and 1878 possible answers. While inadvertently smaller than natural language datasets, further investigation reveals that it covers 88.8% and 70.6% of VQA questions and answers respectively, corroborating its wide diversity. A wide selection of dataset visualizations is provided in section 7.

We associate each question with two types: structural and semantic. The **structural type** is derived from the final operation in the question’s functional program. It can be (1) **verify** for yes/no questions, pertaining to object existence, attribute or relation verification, (2) **query** for all open questions, (3) **choose** for questions that present two alternatives to choose from, e.g. “Is it red or blue?”; (4) **logical** which involve logical inference, and (5) **compare** for comparison questions between two or more objects. The **semantic type** refers to the main subject of the question: (1) *object*: for existence questions, (2) *attribute*: consider the properties or position of an object, (3) *category*: related to object identification within some class, (4) *relation*: for questions asking about the subject or object of a described relation (e.g.

“what is the girl wearing?”), and (5) *global*: about overall properties of the scene such as weather or place. As shown in figure 6, the questions have diverse set of types in both the semantic and structural levels.

We proceed by performing a head-to-head comparison with the VQA2.0 dataset [10], the findings of which are summarized in table 1. Apart from the higher average question length, we can see that GQA consequently contains more verbs and prepositions than VQA (as well as more nouns and adjectives), providing further evidence for its increased compositionality. Semantically, we can see that the GQA questions are significantly more compositional than VQA’s, and involve variety of reasoning skills in much higher frequency (spatial, logical, relational and comparative).

Some VQA question types are not covered by GQA, such as intention (*why*) questions or ones involving OCR or external knowledge. The GQA dataset focuses on factual questions and multi-hop reasoning in particular, rather than covering all types. Comparing to VQA, GQA questions are objective, unambiguous, more compositional and can be answered from the images only, potentially making this benchmark more controlled and convenient for making research progress on.

4.2. Baseline Experiments

We analyse an assortment of models on GQA, including both baselines as well as state-of-the-art models. The baselines include a “blind” LSTM model with access to the questions only, a “deaf” CNN model with access to the images only, an LSTM+CNN model, and two prior models based on the question group, local or global, which return the most common answer for each group, as defined in section 3.4. Beside these, we evaluate the performance of the bottom-up attention model [4] – the winner of 2017 VQA challenge, and the MAC model [11] – a compositional attention state-of-the-art model for CLEVR [15]. For human evaluation, we used Amazon Mechanical Turk to collect human responses for 4000 random questions, taking majority over 5 answers per question. Further description of the evaluated models along with implementation details can be found in section 9.

The evaluation results, including the overall accuracy and the accuracies for each question type, are summarized in table 2. As we can see, the priors and the blind LSTM model achieve very low results of 41.07%: inspection of specific question types reveals that LSTM achieves only 22.7% for open *query* questions, and not far above chance for all other binary question types. We can further see that the “deaf” CNN model achieves as well low results across almost all question types, as expected. On the other hand, state-of-the-art models such as MAC [11] and Bottom-Up Attention [4] perform much better than baselines, but still

Aspect	VQA	GQA
question length	6.2 + 1.9	7.9 + 3.1
verbs	1.4 + 0.6	1.6 + 0.7
nouns	1.9 + 0.9	2.5 + 1.0
adjectives	0.6 + 0.6	0.7 + 0.7
prepositions	0.5 + 0.6	1 + 1
relation questions	19.5%	51.6%
spatial questions	8%	22.4%
logical questions	6%	19%
comparative questions	1%	3%
compositional questions	3%	52%

Table 1: A head-to-head comparison between GQA and VQA. The GQA questions are longer on average, and consequently have more verbs, nouns, adjectives and prepositions than VQA, alluding to their increased compositionality. In addition, GQA demands increased reasoning (spatial, logical, relational and comparative) and includes significantly more compositional questions.

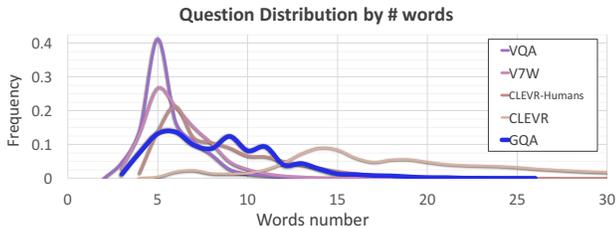


Figure 7: Question length distribution for Visual Question Answering datasets: we can see that GQA questions have a wide range of lengths and are longer on average than all other datasets except the synthetic CLEVR. Note that the long CLEVR questions tend to sound unnatural at times.

well below human scores, offering ample opportunity for further research in the visual reasoning domain.

4.3. Transfer Performance

We tested the transfer performance between the GQA and VQA datasets, training on one and testing on the other: A MAC model trained on GQA achieves 52.1% on VQA before fine-tuning and 60.5% afterwards. Compare these with 51.6% for LSTM+CNN and 68.3% for MAC, when both are trained and tested on VQA. These quite good results demonstrate the realism and diversity of GQA questions, showing the dataset can serve as a good proxy for human-like questions. In contrast, MAC trained on VQA gets 39.8% on GQA before fine-tuning and 46.5% afterwards, illustrating the further challenge GQA poses.

4.4. New Evaluation Metrics

Apart from the standard accuracy metric, and the more detailed type-based diagnosis our dataset supports, we in-

¹¹In a previous version we computed validity, plausibility and consistency in a different way, leading to lower scores than those reported here.

roduce five new metrics to get further insight into visual reasoning methods and point to missing capabilities we believe coherent reasoning models should possess.

Consistency. measure responses consistency across different questions. Recall that in section 3.4, we used the questions’ semantic representation to derive equivalence and entailment relations between them. When being presented with a new question, any learner striving to be trustworthy should not contradict its previous answers. It should not answer “green” to a new question about an apple it has just identified as “red”.

For each question-answer pair (q, a) , we define a set $E_q = q_1, q_2, \dots, q_n$ of entailed questions, the answers to which can be unambiguously inferred given (q, a) . For instance, given the question-answer pair “*Is there a red apple to the left of the white plate?*”, we can infer the answers to questions such as *Is the plate to the right of the apple?*, “*Is there a red fruit to the left of the plate?*”, “*What is the white thing to the right of the apple?*”, etc. For each question q in Q – the set of questions *the model answered correctly*, we measure the model’s accuracy over the entailed questions E_q and then average these score across all questions in Q .

We can see that while people have exceptional consistency of 98.4%, even best models are inconsistent in about 1 out of 5 questions, and models such as LSTM contradict themselves almost half the times. Apparently, achieving high consistency requires deeper understanding of the question semantics in the context of the image, and, in contrast with accuracy, is more robust against “random” educated guesses as it inspects connections between related questions, and thus may serve as a better measure of models’ true visual understanding skills.

Validity and Plausibility. The validity metric checks whether a given answer is in the scope of the question, *e.g.* responding some color to a color question. The plausibility score goes a step further, measuring whether the answer is reasonable, or makes sense, given the question (*e.g.* elephant usually do not eat, say, pizza). Specifically, we check whether the answer occurs at least once in relation with the question’s subject, across the whole dataset, thus, for instance, we consider *e.g.* *red* and *green* as plausible apple colors, whereas *purple* as not.¹² The experiments show that models fail to respond with plausible or even valid answers in at least 5-15% of the times, indicating limited comprehension of some questions. Given that these properties are noticeable statistics of the dataset’s conditional answer distribution, not even depending on the specific images, we would expect a sound method to achieve higher scores.

Distribution. To get further insight into the extent to which methods manage to model the conditional answer

¹²While the plausibility metric may not be fully precise especially for infrequent objects due to potential data scarcity issues, it may provide a good sense of the general level of world-knowledge the model has acquired.

Metric	Global Prior	Local Prior	CNN	LSTM	CNN+LSTM	BottomUp	MAC	Humans
Accuracy	28.93	31.31	17.82	41.07	46.55	49.74	54.06	89.3
Open	16.52	16.99	1.74	22.69	31.80	34.83	38.91	87.4
Binary	42.99	47.53	36.05	61.90	63.26	66.64	71.23	91.2
Query	16.52	16.99	1.55	22.69	31.80	34.83	38.91	87.4
Compare	35.59	41.91	36.34	57.79	56.62	56.32	60.04	93.1
Choose	17.45	26.58	0.85	57.15	61.40	66.56	70.59	94.3
Logical	50.32	50.11	47.18	61.73	62.05	64.03	69.99	88.5
Verify	53.40	58.80	47.02	65.78	67.00	71.45	75.45	90.1
Global	24.70	20.19	8.64	27.22	56.57	60.29	60.82	92.3
Object	49.96	54.00	47.33	74.33	75.90	78.45	81.49	88.1
Attribute	34.89	42.67	22.66	48.28	50.91	53.88	59.82	90.7
Relation	22.88	20.16	11.60	33.24	39.45	42.84	46.16	89.2
Category	15.26	17.31	3.56	22.33	37.49	41.18	44.38	90.3
Distribution	130.86	21.56	19.99	17.93	7.46	5.98	5.34	-
Grounding	-	-	-	-	-	78.47	82.24	-
Validity ¹³	89.02	84.44	35.78	96.39	96.02	96.18	96.16	98.9
Plausibility ¹³	75.34	84.42	34.84	87.30	84.25	84.57	84.48	97.2
Consistency ¹³	51.78	54.34	62.40	68.68	74.57	78.71	81.59	98.4

Table 2: Results for baselines and state-of-the-art models on the GQA dataset. All results refer to the test set. Models are evaluated for overall accuracy as well as accuracy per type. In addition, they are evaluated by validity, plausibility, distribution, consistency, and when possible, grounding metrics. Refer to the text for further detail.

distribution, we define the distribution metric, which measures the overall match between the true answer distribution and the model predicted distribution. For each question global group (section 3.4), we compare the golden and prediction distributions using Chi-Square statistic [21], and then average across all the groups. It allows us to see if the model predicts not only the most common answers but also the less frequent ones. Indeed, the experiments demonstrate that better models such as the state-of-the-art Bottom Up and MAC models score lower than the baselines (for this metric, lower is better), indicating increased capacity in fitting more subtle trends of the dataset’s distribution.

Grounding. For attention-based models, the grounding score checks whether the model attends to regions within the image that are relevant to the question. For each dataset instance, we define a pointer r to the visual region which the question or answer refer to, and compare it to the look at the model visual attention (summing up the overall attention the model gives to r). This metric allows us to evaluate the degree to which the model grounds its reasoning in the image, rather than just making educated guesses based on language priors or world tendencies.

We can see from the experiments that in fact the attention models attend mostly to the right and relevant regions in the image, with grounding scores of about 80%. To verify the reliability of the metric, we further perform experiments with spatial features instead of the object-informed ones used by BottomUp [4] and MAC [11], which lead to a much lower 43% score, demonstrating that indeed object-based features provide models with better granularity for the task, allowing them to focus on more pertinent regions than the coarser spatial features.

5. Conclusion

In this paper, we introduced the GQA dataset for visual reasoning and compositional question answering. We described the dataset generation process, provided baseline experiments and defined new measures to get more insight into models’ behavior and performance. We believe this benchmark can help driving VQA research in the right directions of deeper semantic understanding, sound reasoning, enhanced robustness and improved consistency. A potential avenue towards such goals may involve more intimate integration between visual knowledge extraction and question answering, two flourishing fields that oftentimes have been pursued independently. We strongly hope that GQA will motivate and support the development of more compositional, interpretable and cogent reasoning models, to advance research in scene understanding and visual question answering.

6. Acknowledgments

We wish to thank Justin Johnson for the discussions about the early versions of this work, and Ross Girshick for his inspirational talk at the VQA workshop 2018. We further would like to thank Ranjay Krishna, Eric Cosatto and Alexandru Niculescu-Mizil for the helpful suggestions and comments. Stanford University gratefully acknowledges the generous support of Facebook Inc. as well as the Defense Advanced Research Projects Agency (DARPA) Communicating with Computers (CwC) program under ARO prime contract no. W911NF15-1-0462 for supporting this work.

References

- [1] A. Agrawal, D. Batra, and D. Parikh. Analyzing the behavior of visual question answering models. In *EMNLP*, pages 1955–1960, 2016. [1](#), [2](#)
- [2] A. Agrawal, D. Batra, D. Parikh, and A. Kembhavi. Dont just assume; look and answer: Overcoming priors for visual question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4971–4980, 2018. [1](#), [3](#)
- [3] A. Agrawal, J. Lu, S. Antol, M. Mitchell, C. L. Zitnick, D. Parikh, and D. Batra. VQA: Visual question answering. *International Journal of Computer Vision*, 123(1):4–31, 2017. [1](#), [2](#), [3](#), [6](#)
- [4] P. Anderson, X. He, C. Buehler, D. Teney, M. Johnson, S. Gould, and L. Zhang. Bottom-up and top-down attention for image captioning and VQA. *arXiv preprint arXiv:1707.07998*, 2017. [1](#), [8](#), [10](#), [16](#), [17](#)
- [5] J. Andreas, M. Rohrbach, T. Darrell, and D. Klein. Neural module networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 39–48, 2016. [1](#)
- [6] Y. Attali and M. Bar-Hillel. Guess where: The position of correct answers in multiple-choice test items as a psychometric variable. *Journal of Educational Measurement*, 40(2):109–128, 2003. [3](#)
- [7] A. Das, H. Agrawal, L. Zitnick, D. Parikh, and D. Batra. Human attention in visual question answering: Do humans and deep networks look at the same regions? *Computer Vision and Image Understanding*, 163:90–100, 2017. [1](#), [2](#)
- [8] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 248–255. Ieee, 2009. [16](#)
- [9] A. Fukui, D. H. Park, D. Yang, A. Rohrbach, T. Darrell, and M. Rohrbach. Multimodal compact bilinear pooling for visual question answering and visual grounding. *arXiv preprint arXiv:1606.01847*, 2016. [1](#)
- [10] Y. Goyal, T. Khot, D. Summers-Stay, D. Batra, and D. Parikh. Making the V in VQA matter: Elevating the role of image understanding in visual question answering. In *CVPR*, pages 6325–6334, 2017. [1](#), [2](#), [3](#), [4](#), [8](#)
- [11] D. A. Hudson and C. D. Manning. Compositional attention networks for machine reasoning. *arXiv preprint arXiv:1803.03067*, 2018. [1](#), [8](#), [10](#), [16](#), [17](#)
- [12] G. Inc. Google books ngram corpus. [5](#), [6](#)
- [13] A. Jabri, A. Joulin, and L. van der Maaten. Revisiting visual question answering baselines. In *European conference on computer vision*, pages 727–739. Springer, 2016. [1](#)
- [14] U. Jain, Z. Zhang, and A. G. Schwing. Creativity: Generating diverse questions using variational autoencoders. In *CVPR*, pages 5415–5424, 2017. [4](#)
- [15] J. Johnson, B. Hariharan, L. van der Maaten, L. Fei-Fei, C. L. Zitnick, and R. Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*, pages 1988–1997. IEEE, 2017. [1](#), [2](#), [3](#), [8](#)
- [16] J. Johnson, R. Krishna, M. Stark, L.-J. Li, D. Shamma, M. Bernstein, and L. Fei-Fei. Image retrieval using scene graphs. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3668–3678, 2015. [4](#)
- [17] K. Kafle and C. Kanan. An analysis of visual question answering algorithms. In *Computer Vision (ICCV), 2017 IEEE International Conference on*, pages 1983–1991. IEEE, 2017. [2](#), [3](#), [4](#)
- [18] K. Kafle and C. Kanan. Visual question answering: Datasets, algorithms, and future challenges. *Computer Vision and Image Understanding*, 163:3–20, 2017. [1](#), [2](#), [18](#)
- [19] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. [17](#)
- [20] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma, et al. Visual genome: Connecting language and vision using crowd-sourced dense image annotations. *International Journal of Computer Vision*, 123(1):32–73, 2017. [1](#), [2](#), [3](#), [4](#)
- [21] H. O. Lancaster and E. Seneta. Chi-square distribution. *Encyclopedia of biostatistics*, 2, 2005. [10](#)
- [22] Q. Li, Q. Tao, S. Joty, J. Cai, and J. Luo. Vqa-e: Explaining, elaborating, and enhancing your answers for visual questions. *arXiv preprint arXiv:1803.07464*, 2018. [2](#)
- [23] 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 *European conference on computer vision*, pages 740–755. Springer, 2014. [4](#)
- [24] A. Mahendru, V. Prabhu, A. Mohapatra, D. Batra, and S. Lee. The promise of premise: Harnessing question premises in visual question answering. *arXiv preprint arXiv:1705.00601*, 2017. [3](#), [4](#)
- [25] M. Malinowski and M. Fritz. A multi-world approach to question answering about real-world scenes based on uncertain input. In *Advances in neural information processing systems*, pages 1682–1690, 2014. [3](#), [4](#)
- [26] G. A. Miller. Wordnet: a lexical database for english. *Communications of the ACM*, 38(11):39–41, 1995. [4](#)
- [27] J. Millman, C. H. Bishop, and R. Ebel. An analysis of test-wiseness. *Educational and Psychological Measurement*, 25(3):707–726, 1965. [3](#)
- [28] J. J. Mondak and B. C. Davis. Asked and answered: Knowledge levels when we won’t take ‘don’t know’ for an answer. *Political Behavior*, 23(3):199–224, 2001. [3](#)
- [29] N. Mostafazadeh, I. Misra, J. Devlin, M. Mitchell, X. He, and L. Vanderwende. Generating natural questions about an image. *arXiv preprint arXiv:1603.06059*, 2016. [4](#)
- [30] D. H. Park, L. A. Hendricks, Z. Akata, A. Rohrbach, B. Schiele, T. Darrell, and M. Rohrbach. Multimodal explanations: Justifying decisions and pointing to the evidence. In *31st IEEE Conference on Computer Vision and Pattern Recognition*, 2018. [2](#)
- [31] J. Pennington, R. Socher, and C. Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014. [16](#)
- [32] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In

Advances in neural information processing systems, pages 91–99, 2015. 5, 16, 18

- [33] Y. Rubner, C. Tomasi, and L. J. Guibas. The earth mover’s distance as a metric for image retrieval. *International journal of computer vision*, 40(2):99–121, 2000. 7
- [34] A. Suhr, S. Zhou, I. Zhang, H. Bai, and Y. Artzi. A corpus for reasoning about natural language grounded in photographs. *arXiv preprint arXiv:1811.00491*, 2018. 2
- [35] D. Teney, P. Anderson, X. He, and A. van den Hengel. Tips and tricks for visual question answering: Learnings from the 2017 challenge. *arXiv preprint arXiv:1708.02711*, 2017. 2, 17
- [36] D. Teney, L. Liu, and A. van den Hengel. Graph-structured representations for visual question answering. *arXiv preprint*, 2017. 2
- [37] B. Thomee, D. A. Shamma, G. Friedland, B. Elizalde, K. Ni, D. Poland, D. Borth, and L.-J. Li. Yfcc100m: The new data in multimedia research. *arXiv preprint arXiv:1503.01817*, 2015. 4
- [38] A. Torralba and A. A. Efros. Unbiased look at dataset bias. In *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pages 1521–1528. IEEE, 2011. 2
- [39] D. Xu, Y. Zhu, C. B. Choy, and L. Fei-Fei. Scene graph generation by iterative message passing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, volume 2, 2017. 4
- [40] Z. Yang, X. He, J. Gao, L. Deng, and A. Smola. Stacked attention networks for image question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 21–29, 2016. 1
- [41] P. Zhang, Y. Goyal, D. Summers-Stay, D. Batra, and D. Parikh. Yin and yang: Balancing and answering binary visual questions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5014–5022, 2016. 1, 3
- [42] S. Zhang, L. Qu, S. You, Z. Yang, and J. Zhang. Automatic generation of grounded visual questions. *arXiv preprint arXiv:1612.06530*, 2016. 4
- [43] Y. Zhu, O. Groth, M. Bernstein, and L. Fei-Fei. Visual7w: Grounded question answering in images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4995–5004, 2016. 1, 2, 3

7. Dataset Visualizations

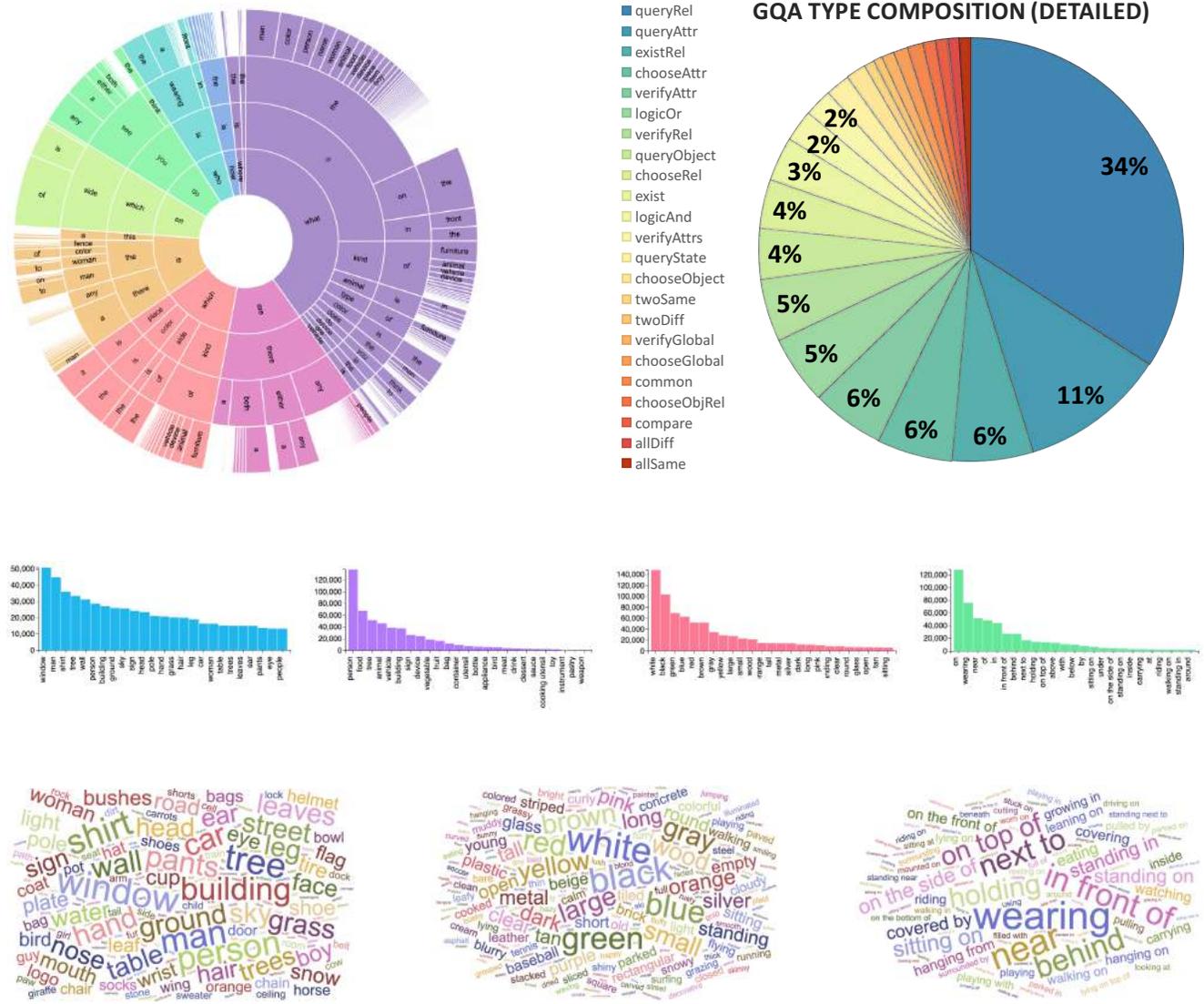


Figure 8: **Top left:** Distribution of GQA questions by first four words. The arc length is proportional to the number of questions containing that prefix. White areas correspond to marginal contributions. **Top right:** question type distribution; please refer to the table 3 for details about each type. **Middle rows:** Occurrences number of the most frequent objects, categories, attributes and relations (excluding left/right). **Third row:** Word clouds for frequent objects, attributes and relations.



GQA

1. What is the **woman** to the right of the **boat** holding? umbrella
2. Are there **men** to the left of the **person** that is holding the **umbrella**? no
3. What color is the **umbrella** the **woman** is holding? purple

VQA

1. Why is the person using an umbrella?
2. Is the picture edited?
3. What's the color of the umbrella?



GQA

1. Is that a **giraffe** or an **elephant**? giraffe
2. Who is feeding the **giraffe** behind the **man**? lady
3. Is there any **fence** near the **animal** behind the **man**? yes
4. On which side of the image is the **man**? right
5. Is the **giraffe** is behind the **man**? yes

VQA

1. What animal is the lady feeding?
2. Is it raining?
3. Is the man wearing sunglasses?



GQA

1. Is the **person's hair** brown and long? yes
2. What **appliance** is to the left of the **man**? refrigerator
3. Is the **man** to the left or to the right of a **refrigerator**? right
4. Who is in front of the **appliance** on the left? man
5. Is there a **necktie** in the picture that is not red? yes
6. What is the **person** in front of the **refrigerator** wearing? suit
7. What is hanging on the **wall**? picture
8. Does the **vest** have different color than the **tie**? no
9. What is the color of the **shirt**? white
10. Is the color of the **vest** different than the **shirt**? yes

VQA

1. Does this man need a haircut?
2. What color is the guys tie?
3. What is different about the man's suit that shows this is for a special occasion?



GQA

1. Who wears the **gloves**? player
2. Are there any **horses** to the left of the **man**? no
3. Is the **man** to the right of the **player** that wears gloves? no
4. Is there a **bag** in the picture? no
5. Do the **hat** and the **plate** have different colors? yes

VQA

1. What is the man holding?
2. Where are the people playing?
3. Is the player safe?
4. What is the sport being played?



GQA

1. What is the **person** doing? playing
2. Is the **entertainment center** at the bottom or at the top? bottom
3. Is the **entertainment center** wooden and small? yes
4. Are the pants blue? no
5. Do you think the **controller** is red? no

VQA

1. What colors are the walls?
2. What game is the man playing?
3. Why do they stand to play?



GQA

1. Are there any **coats**? yes
2. Do you see a red **coat** in the image? no
3. Is the **person** that is to the left of the **man** exiting a **truck**? no
4. Which place is this? road

VQA

1. Where is the bus driver?
2. Why is the man in front of the bus?
3. What numbers are repeated in the bus number?



GQA

1. What is in front of the green **fence**? gate
2. Of which color is the **gate**? silver
3. Where is this? street
4. What color is the **fence** behind the **gate**? green
5. Is the **fence** behind the **gate** both brown and metallic? no

VQA

1. What are the yellow lines called?
2. Why don't the trees have leaves?
3. Where is the stop sign?

Figure 10: Examples of questions from GQA and VQA, for the same images. As the examples demonstrate, GQA questions tend to involve more elements from the image compared to VQA questions, and are longer and more compositional as well. Conversely, VQA questions tend to be a bit more ambiguous and subjective, at times with no clear and conclusive answer. Finally, we can see that GQA provides more questions for each image and thus covers it more thoroughly than VQA.

8. Dataset Balancing

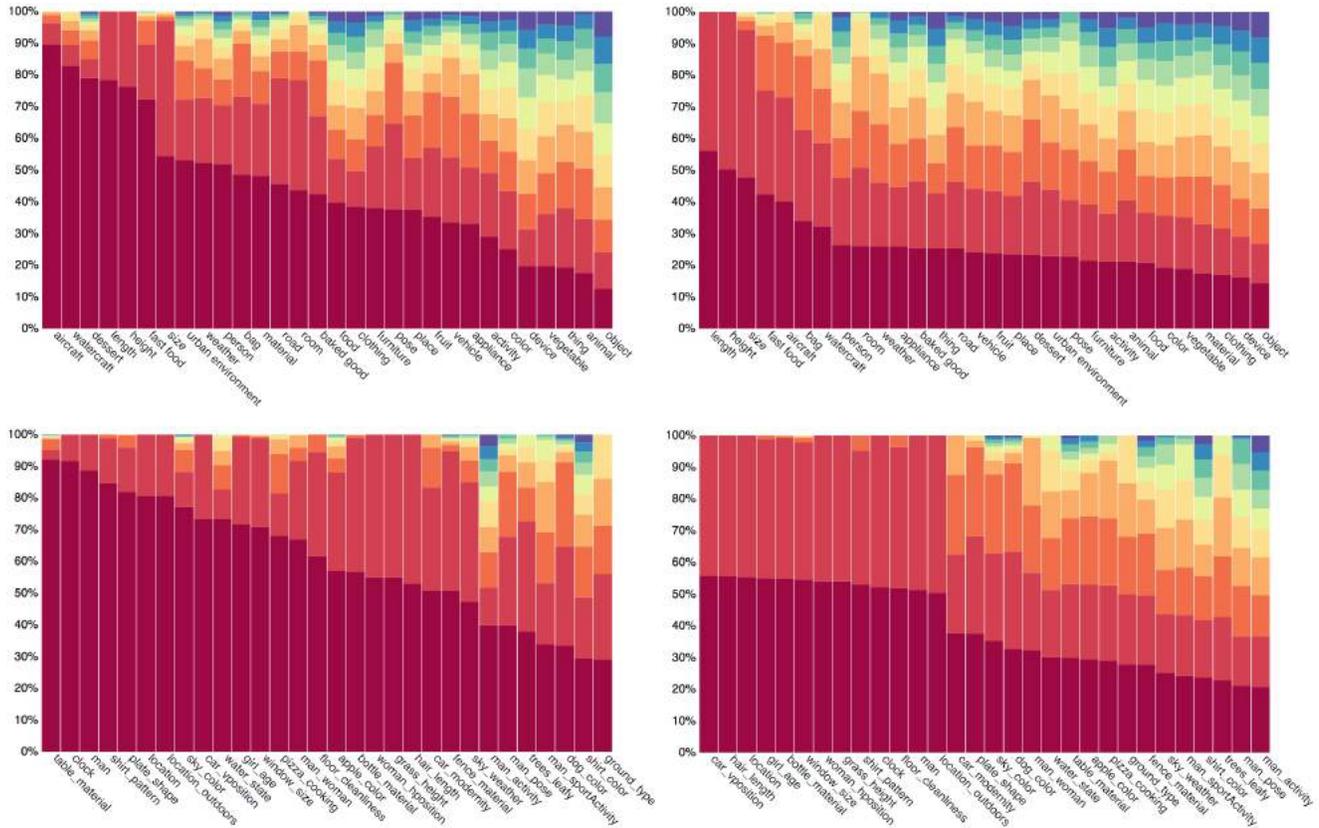


Figure 11: Impact of the dataset balancing on the conditional answer distribution: The left side shows the distribution before any balancing. We show the top 10 answers for a selection of question groups, where the column height corresponds to the relative frequency of each answer. The top row shows global question groups such as color questions, questions about animals, etc. while the bottom row shows local ones *e.g.* *apple-color*, *table-material* etc (section 3.4). Indeed, we can see that the distributions are heavily biased. The right side shows the distributions after balancing, more uniform and with heavier tails, while intentionally retaining the original real-world tendencies up to a tunable degree.

As discussed in section 3.5, given the original 22M auto-generated questions, we have performed answer-distribution balancing, similarities reduction and type-based sampling, reducing its size to a 1.7M balanced dataset. The balancing is performed in an iterative manner: as explained in section 3.4, for each question group (*e.g.* color questions), we iterate over the answer distribution, from the most to least frequent answers: (a_i, c_i) when a_i is the answer and c_i is its count. In each iteration i , we downsample the head distribution $(a_j, j \leq i)$ such that the ratio between the head and its complementary tail will be bounded by b . While doing so, we also make sure to set minimum and maximum bounds on the frequency ratio $\frac{c_{i+1}}{c_i}$ of each pair of consequent answers a_i, a_{i+1} . The results of this process is shown in figure 11. Indeed we can see how the distribu-

tion is “pushed” away from the head and spreads over the tail, while intentionally maintaining the original real-world tendencies presented in the data, to retain its authenticity.

9. Baselines Implementation Details

In section 4.2, we perform experiments over multiple baselines and state-of-the-art models. All CNN models use spatial features pre-trained on ImageNet [8], whereas state-of-the-art approaches such as bottomUp [4] and MAC [11] are based on object-based features produced by faster R-CNN detector [32]. All models use GloVe word embeddings of dimension 300 [31]. To allow a fair comparison, all the models use the same LSTM, CNN and classifier components, and so the only difference between the models stem from their core architectural design.

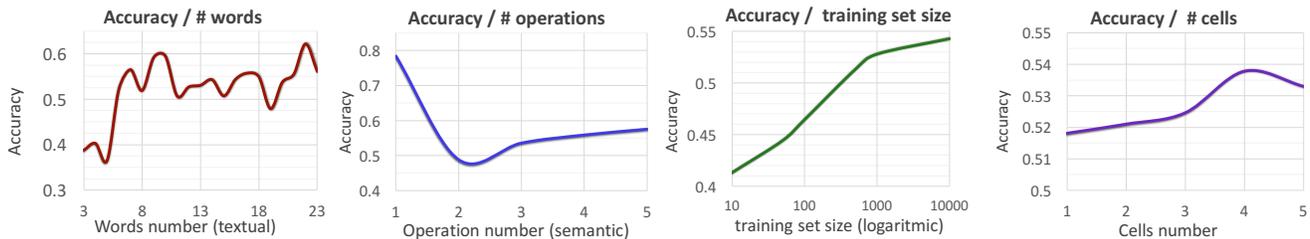


Figure 12: From left to right: (1) Accuracy as a function of textual question length – the number of words in the question. (2) Accuracy as a function of semantic question length – the number of operations in its functional program. (3) Performance as a function of the subset size used for training, ranging from 10K to 10M. (4) Accuracy for different lengths of MAC networks, suggesting that indeed GQA questions are compositional.

We have used a sigmoid-based classifier and trained all models using Adam [19] for 15 epochs, each takes about an hour to complete. For MAC [11], we use the official authored code available online, with 4 cells. For BottomUp [4], since the official implementation is unfortunately not publicly available, we re-implemented the model, carefully following details presented in [4, 35]. To ensure the correctness of our implementation, we have tested the model on the standard VQA dataset, achieving 67%, which matches the original scores reported by Anderson *et al.* [4].

10. Further Diagnosis

Following section 4.2, and in order to get more insight into models’ behaviors and tendencies, we perform further analysis of the top-scoring model for the GQA dataset, MAC [11]. The MAC network is a recurrent attention network that reasons in multiple concurrent steps over both the question and the image, and is thus geared towards compositional reasoning as well as rich scenes with several regions of relevance.

We assess the model along multiple axes of variation, including question length, both textually, *i.e.* number of words, and semantically, *i.e.* number of reasoning operations required to answer it, where an operation can be *e.g.* following a relation from one object to another, attribute identification, or a logical operation such as *or*, *and* or *not*. We provide additional results for different network lengths (namely, cells number) and varying training-set sizes, all can be found in figure 12.

Interestingly, question textual length correlates positively with the model accuracy. It may be the case that longer questions reveal more cues or information that the model can exploit, potentially sidestepping direct reasoning about the image. However, question semantic length has the opposite impact as expected: 1-step questions are particularly easy for models than the compositional ones which involve more steps.

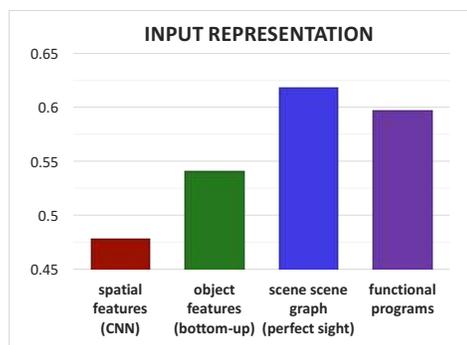


Figure 13: Performance as a function of the input representation. We encode the scenes through three different methods: spatial features produced by a standard pretrained CNN, object-based features generated by a faster R-CNN detector, and direct embedding of the scene graph semantic representation, equivalent to having perfect sight. We further experiment with both textual questions as well as their counterpart functional programs as input. We can see that the more semantically-imbued the representations get, the higher the accuracy obtained.

GQA SEMANTIC STEPS

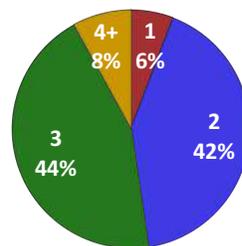


Figure 14: Distribution of GQA questions semantic length (number of computation steps to arrive at the answer). We can see that most questions require about 2-3 reasoning steps, where each step may involve tracking a relation between objects, an attribute identification or a logical operation.

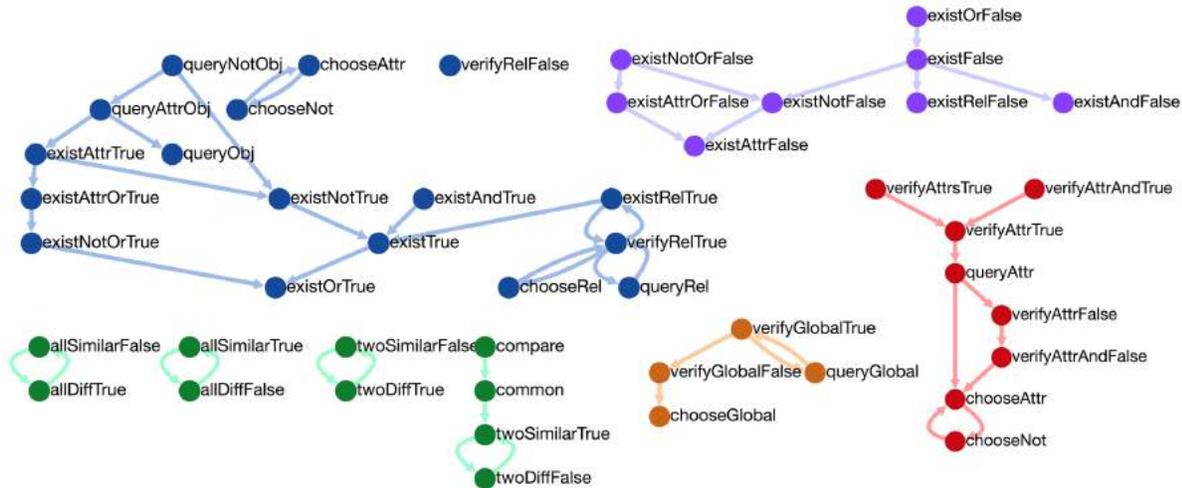


Figure 15: Entailment relations between different question types. In section 3.4 we discuss the entailment and equivalences between questions. Since every question in the dataset has a matching logical representation of the sequence of reasoning steps, we can formally compute all the entailment and equivalence relations between different questions. Indeed, a cogent and reasonable learner should be consistent between its own answers, *e.g.* should not answer “red” to a question about the color of an object it has just identified as blue. Some more subtle relations also occur, such as those involving relations, *e.g.* if X is above Y, then Y is below X, and X is not below Y, etc. figure 15 shows all the logical relations between the various question types. Refer to table 3 for a complete catalog of the different types. Experiments show that while people excel at consistency, achieving the impressive 98.4%, deep learning models perform much worse in this task, with 69% - 82%. These results cast a doubt about the reliability of existing models and their true visual understanding skills. We therefore believe that improving their skills towards enhanced consistency and cogency is an important direction, which we hope our dataset will encourage.

We can further see that longer MAC networks with more cells are more competent in performing the GQA task, substantiating its increased compositionality. Other experiments show that increasing the training set size has significant impact on the model’s performance, as found out also by Kafle *et al.* [18]. Apparently, the training set size has not reached saturation yet and so models may benefit from even larger datasets.

Finally, we have measured the impact of different input representations on the performance. We encode the visual scene with three different methods, ranging from standard pretrained CNN-based spatial features, to object-informed features obtained through faster R-CNNs detectors [32], up to even a “perfect sight” model that has access to the precise semantic scene graph through direct node and edge embeddings. As figure 12 shows, the more high-level and semantic the representation is, the better are the results.

On the question side, we explore both training on the standard textual questions as well as the semantic functional programs. MAC achieves 53.8% accuracy and 81.59% consistency on the textual questions and 59.7% and 85.85% on the programs, demonstrating the usefulness and further challenge embodied in the former. It is also more consistent. Indeed, the programs consist of only a small operations vocabulary, whereas the questions use both synonyms and hundreds of possible structures, incorporating probabilistic rules to make them more natural and diverse. In particular, GQA questions have sundry subtle and challenging linguistic phenomena such as long-range dependencies, absent from the canonical programs. The textual questions thus provide us with the opportunity to engage with real, interesting and significant aspects of natural language, and consequently foster the development of models with enhanced language comprehension skills.

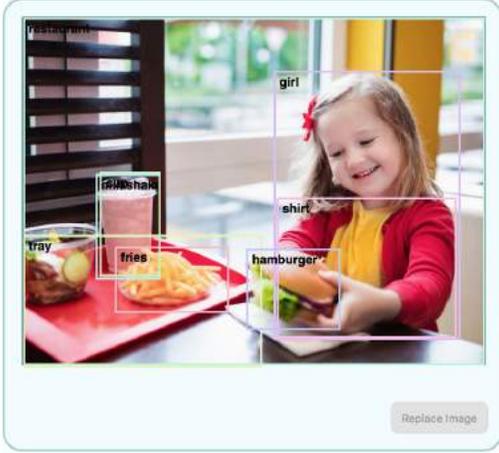
Image Annotation

In this HIT, you are going to annotate objects in 4 images, list the properties of each object, and the relations between them. You will do it in multiple steps:

1. Draw a box around each object in the image on the left, and type its name. Please mark as many objects in the image as possible (usually ~12-20 objects), and make sure to create a tight box around each object, the smallest to cover it. Objects can be people and animals, foods and drinks, clothing items, furniture, appliances, vehicles, buildings, places etc.
2. Write up to five properties for each object on the right side. Properties can be any adjectives, colors, materials, size, length, shape, activity (running, sleeping), etc.
3. After clicking next, write relations between pairs of objects in the image, for instance "girl, eating, cake". Each relation (eating) goes from a source object (girl), to a target object (cake). Relations can be verbs (holding, chasing) spatial relations or prepositions (on top of, around, behind). Please try to write as many relations as possible (usually 10-15 relations), but it's ok if you can't find enough.

*If there's a group of close same-type objects you should mark them together (e.g. "fries").

*In case the word you type is not in our vocabulary, a few similar alternatives will be presented to select from. Bonuses for careful work! (bad work may be rejected)



Here is an example of a few annotated objects, along with their properties:

Object	Properties
1. restaurant	modern, clean, bright
2. milkshake	pink, bright, sweet
3. girl	young, blond, happy, sitting
4. shirt	red, cloth, long sleeved
5. tray	red, plastic, rectangular
6. fries	yellow, cooked, thin
7. cup	plastic, large, transparent, full
8. hamburger	leasty

Other objects that have to be annotated are: window, table, napkin, salad and bowl.

After annotating all the objects, you will have to write relations between them:

1. girl	holding	hamburger
2. fries	on top of	tray
3. girl	wearing	shirt
4. cup	contain	milkshake

Bonuses will be given for good work, with many objects, properties and relations! (but bad work may be rejected).

Previous Next

Image Question Answering

In this HIT, you are going to answer questions about pictures!

We will show you 4 pictures and 5-10 questions about each of them (the same picture may appear twice).

- For each question, start by typing your answer in the text box right to it. If you don't know the answer, please type "I don't know".
- In case your answer is not one of the possible answers in our system, a few relevant alternatives will be shown to choose from. Please select the one that sounds the most correct to you among them. If you believe none of the choices is right please select "None of the above".
- The answers are usually short, about 1-2 words.

P.S. You'll receive bonus for each question you answer correctly! So try to do your best! :) Good Luck!



1. Is there any milk in the bowl left of the apple?
2. Is the bowl right of the green apple?
3. Are there red apples in this picture?
4. Which color do you think is the apple?
5. What type of fruit in the image is round?
6. What color is the fruit on the right side, re
7. On which side of the photo is the apple, th
8. Is there a spoon right of the food in the ce
9. Which color do you think is the apple?
10. Are there red apples in this picture?

1 / 4 Previous Next

Figure 16: The interfaces used for human experiments on Amazon Mechanical Turk. **Top:** Each HIT displays several images and asks turkers to list objects and annotate their corresponding bounding boxes. In addition, the turkers are requested to specify attributes and relations between the objects. An option to switch between images is also given to allow the turkers to choose rich enough images to work on. **Bottom:** Each HIT displays multiple questions and requires the turkers to respond. Since there is a closed set of possible answers (from a vocabulary with Approximately 1878 tokens), and in order to allow a fair comparison between human and models' performance, we give turkers the option to respond in unconstrained free-form language, but also suggest them multiple answers from our vocabulary that are the most similar to theirs (using word embedding distances). However, turkers are not limited to choose from the suggestions in case they believe none of the proposed answers is correct.