


# Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations

Ranjay Krishna<sup>1</sup>  · Yuke Zhu<sup>1</sup> · Oliver Groth<sup>2</sup> · Justin Johnson<sup>1</sup> · Kenji Hata<sup>1</sup> · Joshua Kravitz<sup>1</sup> · Stephanie Chen<sup>1</sup> · Yannis Kalantidis<sup>3</sup> · Li-Jia Li<sup>4</sup> · David A. Shamma<sup>5</sup> · Michael S. Bernstein<sup>1</sup> · Li Fei-Fei<sup>1</sup>

Received: 23 February 2016 / Accepted: 12 September 2016 / Published online: 6 February 2017  
© The Author(s) 2017. This article is published with open access at Springerlink.com

**Abstract** Despite progress in perceptual tasks such as image classification, computers still perform poorly on cognitive tasks such as image description and question answering. Cognition is core to tasks that involve not just recognizing, but reasoning about our visual world. However, models used to tackle the rich content in images for cognitive tasks are still being trained using the same datasets designed for perceptual tasks. To achieve success at cognitive tasks, models need to understand the interactions and relationships between objects in an image. When asked “What vehicle is the person riding?”, computers will need to identify the objects in an image as well as the relationships *riding(man, carriage)* and *pulling(horse, carriage)* to answer correctly that “the person is riding a horse-drawn carriage.” In this paper, we present the Visual Genome dataset to enable the modeling of such relationships. We collect dense annotations of objects, attributes, and relationships within each image to learn these models. Specifically, our dataset contains over 108K images where each image has an average of 35 objects, 26 attributes, and 21 pairwise relationships between objects. We canonicalize the objects, attributes, relationships, and noun phrases in region descriptions and questions answer pairs to WordNet synsets. Together, these annotations represent the densest

and largest dataset of image descriptions, objects, attributes, relationships, and question answer pairs.

**Keywords** Computer vision · Dataset · Image · Scene graph · Question answering · Objects · Attributes · Relationships · Knowledge · Language · Crowdsourcing

## 1 Introduction

A holy grail of computer vision is the complete understanding of visual scenes: a model that is able to name and detect objects, describe their attributes, and recognize their relationships. Understanding scenes would enable important applications such as image search, question answering, and robotic interactions. Much progress has been made in recent years towards this goal, including image classification (Perronnin et al. 2010; Simonyan and Zisserman 2014; Krizhevsky et al. 2012; Szegedy et al. 2015) and object detection (Girshick et al. 2014; Sermanet et al. 2013; Girshick 2015; Ren et al. 2015b). An important contributing factor is the availability of a large amount of data that drives the statistical models that underpin today’s advances in computational visual understanding. While the progress is exciting, we are still far from reaching the goal of comprehensive scene understanding. As Fig. 1 shows, existing models would be able to detect discrete objects in a photo but would not be able to explain their interactions or the relationships between them. Such explanations tend to be *cognitive* in nature, integrating *perceptual* information into conclusions about the relationships between objects in a scene (Bruner 1990; Firestone and Scholl 2015). A cognitive understanding of our visual world thus requires that we complement computers’ ability to detect objects with abilities to describe those

Communicated by Margaret Mitchell, John Platt, and Kate Saenko.

✉ Ranjay Krishna  
ranjaykrishna@cs.stanford.edu

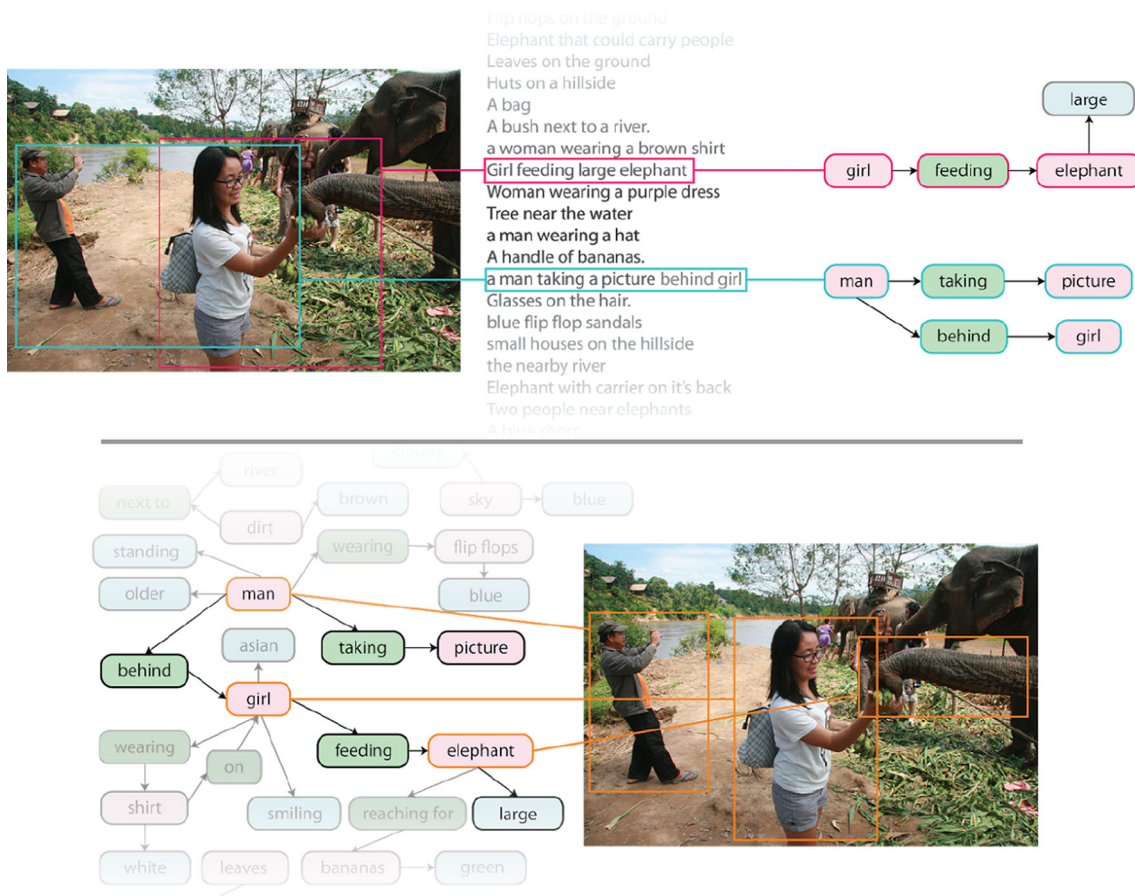
<sup>1</sup> Stanford University, Stanford, CA, USA

<sup>2</sup> Dresden University of Technology, Dresden, Germany

<sup>3</sup> Yahoo Inc., San Francisco, CA, USA

<sup>4</sup> Snapchat Inc., Los Angeles, CA, USA

<sup>5</sup> Centrum Wiskunde & Informatica (CWI), Amsterdam, The Netherlands



**Fig. 1** An overview of the data needed to move from perceptual awareness to cognitive understanding of images. We present a dataset of images densely annotated with numerous region descriptions, objects, attributes, and relationships. Some examples of region descriptions (e.g.

“girl feeding large elephant” and “a man taking a picture behind girl”) are shown (top). The objects (e.g. elephant), attributes (e.g. large) and relationships (e.g. feeding) are shown (bottom). Our dataset also contains image related question answer pairs (not shown)

objects (Isola et al. 2015) and understand their interactions within a scene (Sadeghi and Farhadi 2011).

There is an increasing effort to put together the next generation of datasets to serve as training and benchmarking datasets for these deeper, cognitive scene understanding and reasoning tasks, the most notable being MS-COCO (Lin et al. 2014) and VQA (Antol et al. 2015). The MS-COCO dataset consists of 300K real-world photos collected from Flickr. For each image, there is pixel-level segmentation of 80 object classes (when present) and 5 independent, user-generated sentences describing the scene. VQA adds to this a set of 614K question answer pairs related to the visual contents of each image (see more details in Sect. 3.1). With this information, MS-COCO and VQA provide a fertile training and testing ground for models aimed at tasks for accurate object detection, segmentation, and summary-level image captioning (Kiros et al. 2014; Mao et al. 2014; Karpathy and Fei-Fei 2015) as well as basic QA (Ren et al. 2015a; Malinowski et al. 2015; Gao et al. 2015; Malinowski and Fritz 2014). For example, a state-of-the-art model (Karpathy and Fei-Fei 2015) provides a description of one MS-COCO image in

Fig. 1 as “two men are standing next to an elephant.” But what is missing is the further understanding of where each object is, what each person is doing, what the relationship between the person and elephant is, etc. Without such relationships, these models fail to differentiate this image from other images of people next to elephants.

To understand images thoroughly, we believe three key elements need to be added to existing datasets: a **grounding of visual concepts to language** (Kiros et al. 2014), a more **complete set of descriptions and QAs** for each image based on multiple image regions (Johnson et al. 2015), and a **formalized representation** of the components of an image (Hayes 1978). In the spirit of mapping out this complete information of the visual world, we introduce the Visual Genome dataset. The first release of the Visual Genome dataset uses 108,077 images from the intersection of the YFCC100M (Thomee et al. 2016) and MS-COCO (Lin et al. 2014). Section 5 provides a more detailed description of the dataset. We highlight below the motivation and contributions of the three key elements that set Visual Genome apart from existing datasets.

The Visual Genome dataset regards relationships and attributes as first-class citizens of the annotation space, in addition to the traditional focus on objects. Recognition of relationships and attributes is an important part of the complete understanding of the visual scene, and in many cases, these elements are key to the story of a scene (e.g., the difference between “a dog chasing a man” versus “a man chasing a dog”). The Visual Genome dataset is among the first to provide a detailed labeling of object interactions and attributes, **grounding visual concepts to language**.<sup>1</sup>

An image is often a rich scenery that cannot be fully described in one summarizing sentence. The scene in Fig. 1 contains multiple “stories”: “a man taking a photo of elephants,” “a woman feeding an elephant,” “a river in the background of lush grounds,” etc. Existing datasets such as Flickr 30K (Young et al. 2014) and MS-COCO (Lin et al. 2014) focus on high-level descriptions of an image.<sup>2</sup> Instead, for each image in the Visual Genome dataset, we collect more than 50 descriptions for different regions in the image, providing a much denser and more **complete set of descriptions of the scene**. In addition, inspired by VQA (Antol et al. 2015), we also collect an average of 17 question answer pairs based on the descriptions for each image. Region-based question answers can be used to jointly develop NLP and vision models that can answer questions from either the description or the image, or both of them.

With a set of dense descriptions of an image and the explicit correspondences between visual pixels (i.e. bounding boxes of objects) and textual descriptors (i.e. relationships, attributes), the Visual Genome dataset is poised to be the first image dataset that is capable of providing a structured **formalized representation** of an image, in the form that is widely used in knowledge base representations in NLP (Zhou et al. 2007; GuoDong et al. 2005; Culotta and Sorensen 2004; Socher et al. 2012). For example, in Fig. 1, we can formally express the relationship *holding* between the woman and food as *holding(woman, food)*. Putting together all the objects and relations in a scene, we can represent each image as a scene graph (Johnson et al. 2015). The scene graph representation has been shown to improve semantic image retrieval (Johnson et al. 2015; Schuster et al. 2015) and image captioning (Farhadi et al. 2009; Chang et al. 2014; Gupta and Davis 2008). Furthermore, all objects, attributes and relationships in each image in the Visual Genome dataset are canonicalized to its corresponding WordNet (Miller 1995) ID (called a synset ID). This mapping connects all images in Visual Genome and provides an effective way to consistently

query the same concept (object, attribute, or relationship) in the dataset. It can also potentially help train models that can learn from contextual information from multiple images (Figs. 2, 3).

In this paper, we introduce the Visual Genome dataset with the aim of training and benchmarking the next generation of computer models for comprehensive scene understanding. The paper proceeds as follows: In Sect. 2, we provide a detailed description of each component of the dataset. Section 3 provides a literature review of related datasets as well as related recognition tasks. Section 4 discusses the crowdsourcing strategies we deployed in the ongoing effort of collecting this dataset. Section 5 is a collection of data analysis statistics, showcasing the key properties of the Visual Genome dataset. Last but not least, Sect. 6 provides a set of experimental results that use Visual Genome as a benchmark.

Further visualizations, API, and additional information on the Visual Genome dataset can be found online.<sup>3</sup>

## 2 Visual Genome Data Representation

The Visual Genome dataset consists of seven main components: *region descriptions*, *objects*, *attributes*, *relationships*, *region graphs*, *scene graphs*, and *question answer pairs*. Figure 4 shows examples of each component for one image. To enable research on comprehensive understanding of images, we begin by collecting descriptions and question answers. These are raw texts without any restrictions on length or vocabulary. Next, we extract objects, attributes and relationships from our descriptions. Together, objects, attributes and relationships comprise our scene graphs that represent a formal representation of an image. In this section, we break down Fig. 4 and explain each of the seven components. In Sect. 4, we will describe in more detail how data from each component is collected through a crowdsourcing platform.

### 2.1 Multiple Regions and Their Descriptions

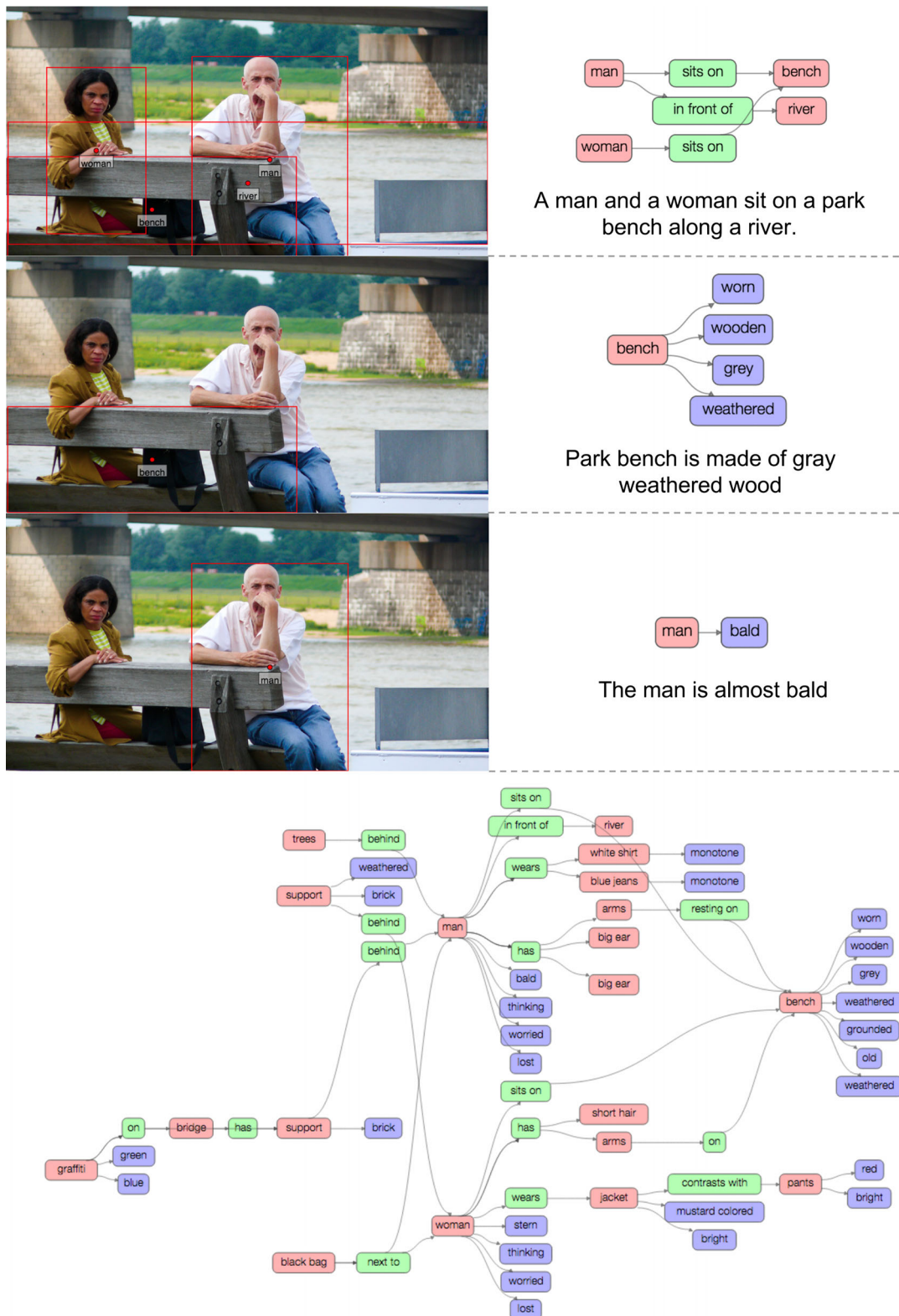
In a real-world image, one simple summary sentence is often insufficient to describe all the contents of and interactions in an image. Instead, one natural way to extend this might be a collection of descriptions based on different regions of a scene. In Visual Genome, we collect diverse human-generated image region descriptions, with each region localized by a bounding box. In Fig. 5, we show three examples of region descriptions. Regions are allowed to have a high degree of overlap with each other when the descriptions differ. For example, “yellow fire hydrant” and “woman in shorts is standing behind the man” have very little

<sup>1</sup> The Lotus Hill Dataset (Yao et al. 2007) also provides a similar annotation of object relationships, see Sec 3.1.

<sup>2</sup> COCO has multiple sentences generated independently by different users, all focusing on providing an overall, one sentence description of the scene.

<sup>3</sup> <https://visualgenome.org>.

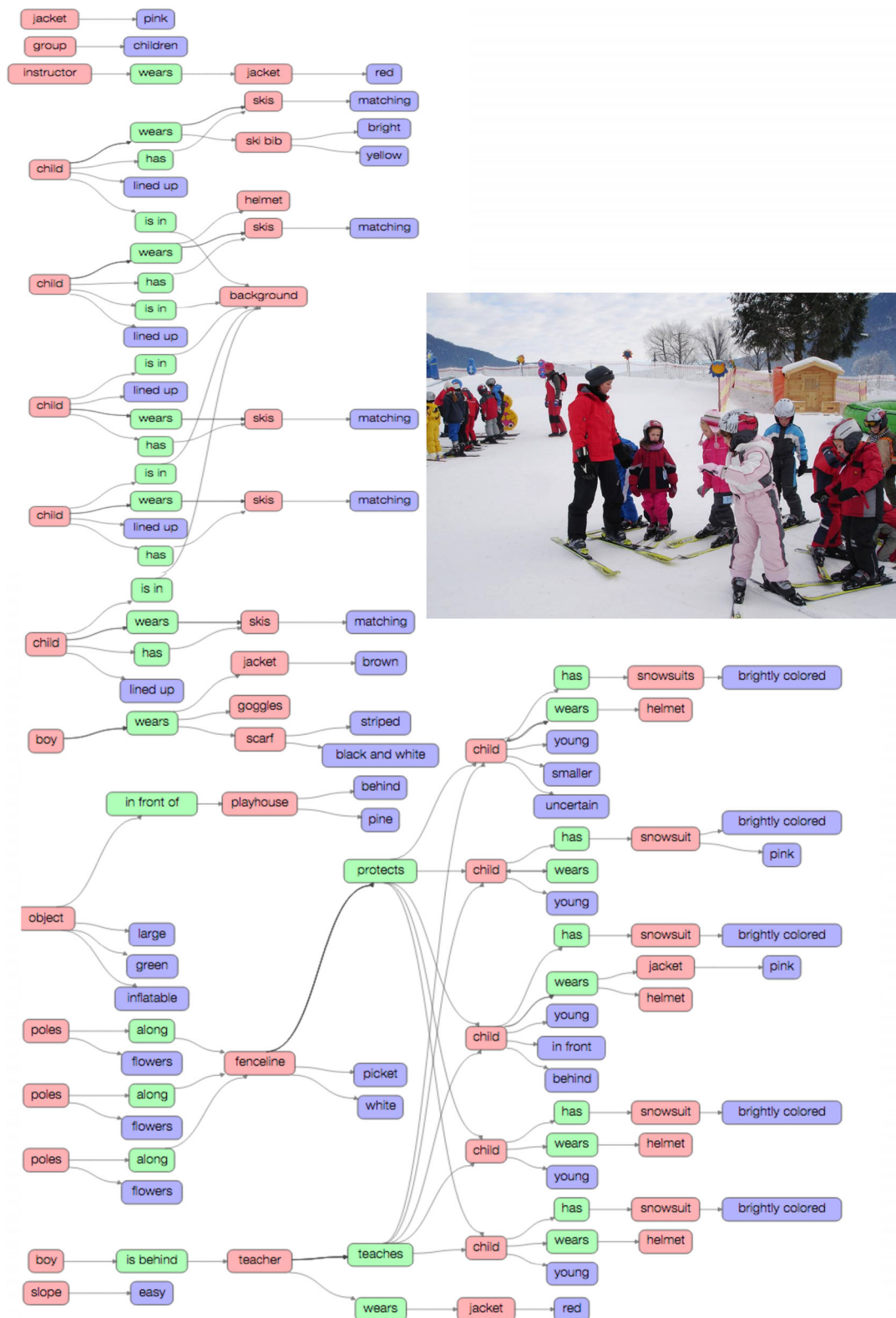




**Fig. 2** An example image from the Visual Genome dataset. We show 3 region descriptions and their corresponding region graphs. We also show the connected scene graph collected by combining all of the image's region graphs. The *top region* description is “a man and a

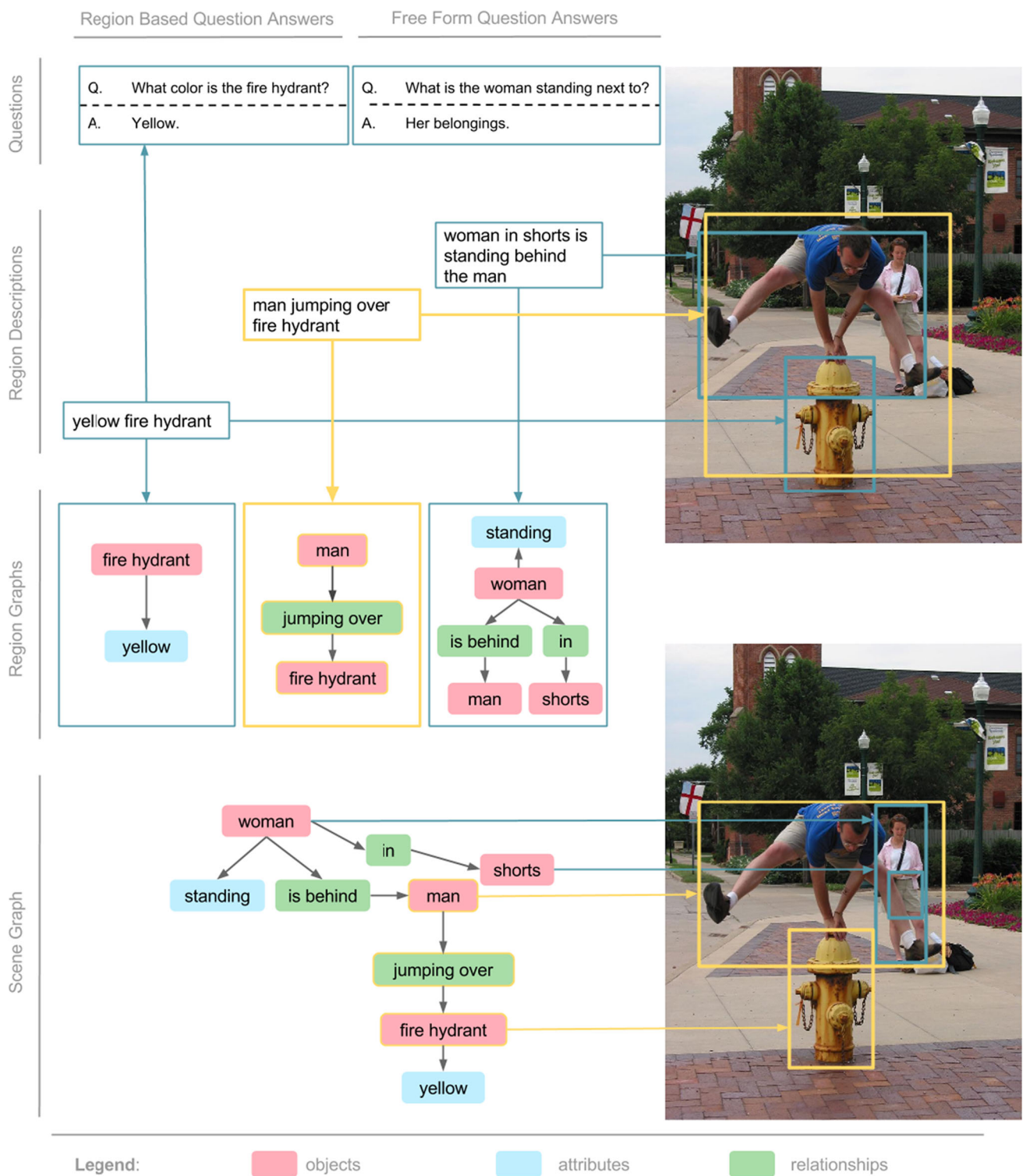
woman sit on a park bench along a river.” It contains the objects: man, woman, bench and river. The relationships that connect these objects are: *sits\_on*(man, bench), *in\_front\_of*(man, river), and *sits\_on*(woman, bench)





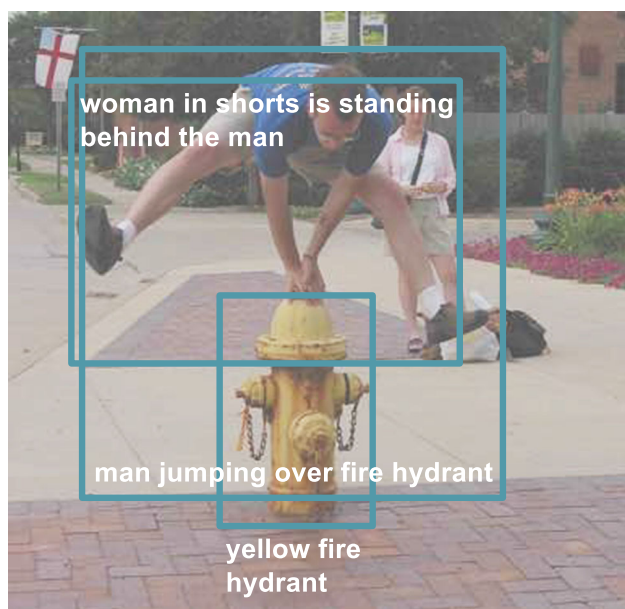
**Fig. 3** An example image from our dataset along with its scene graph representation. The scene graph contains objects (child, instructor, helmet, etc.) that are localized in the image as bounding boxes (not shown). These objects also have attributes: large,

green, behind, etc. Finally, objects are connected to each other through relationships: *wears(child, helmet)*, *wears(instructor, jacket)*, etc



**Fig. 4** A representation of the Visual Genome dataset. Each image contains region descriptions that describe a localized portion of the image. We collect two types of question answer pairs (QAs): freeform QAs and region-based QAs. Each region is converted to a region graph

representation of objects, attributes, and pairwise relationships. Finally, each of these region graphs are combined to form a scene graph with all the objects grounded to the image. *Best viewed in color*



**Fig. 5** To describe all the contents of and interactions in an image, the Visual Genome dataset includes multiple human-generated image regions descriptions, with each region localized by a bounding box. Here, we show three regions descriptions on various image regions: “man jumping over a fire hydrant,” “yellow fire hydrant,” and “woman in shorts is standing *behind* the man”

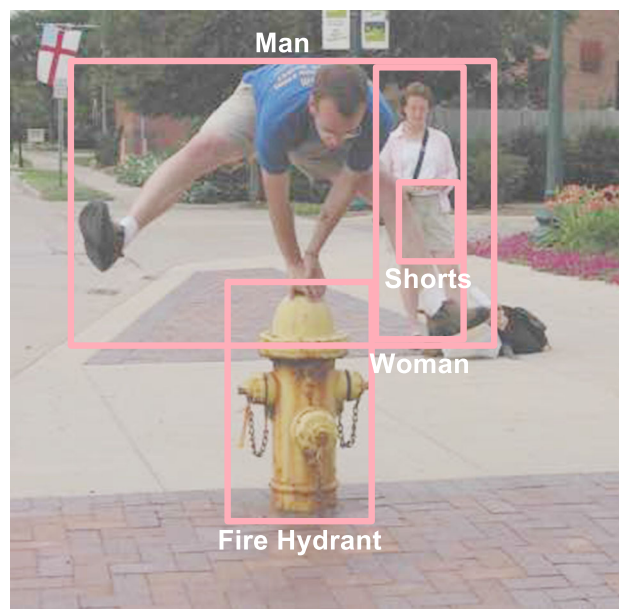
overlap, while “man jumping over fire hydrant” has a very high overlap with the other two regions. Our dataset contains on average a total of 50 region descriptions per image. Each description is a phrase ranging from 1 to 16 words in length describing that region.

## 2.2 Multiple Objects and Their Bounding Boxes

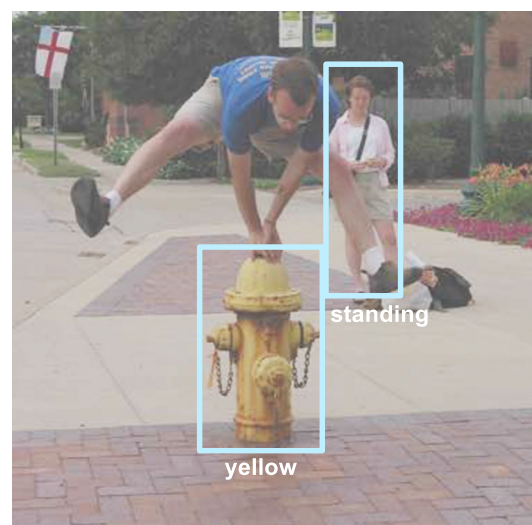
Each image in our dataset consists of an average of 35 objects, each delineated by a tight bounding box (Fig. 6). Furthermore, each object is canonicalized to a synset ID in WordNet (Miller 1995). For example, man would get mapped to `man.n.03` (the generic use of the word to refer to any human being). Similarly, person gets mapped to `person.n.01` (a human being). Afterwards, these two concepts can be joined to `person.n.01` since this is a hypernym of `man.n.03`. We did not standardize synsets in our dataset. However, given our canonicalization, this is easily possible leveraging the WordNet ontology to avoid multiple names for one object (e.g. man, person, human), and to connect information across images.

## 2.3 A Set of Attributes

Each image in Visual Genome has an average of 26 attributes. Objects can have zero or more attributes asso-



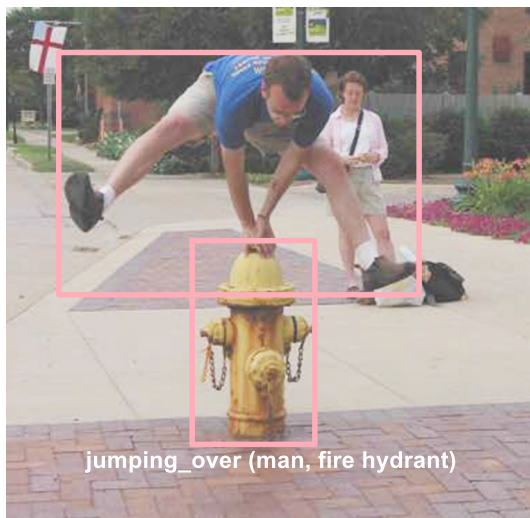
**Fig. 6** From all of the region descriptions, we extract all objects mentioned. For example, from the region description “man jumping *over* a fire hydrant,” we extract man and fire hydrant



**Fig. 7** Some descriptions also provide attributes for objects. For example, the region description “yellow fire hydrant” adds that the fire hydrant is yellow. Here we show two attributes: yellow and standing

ciated with them. Attributes can be color (e.g. yellow), states (e.g. standing), etc. (Fig. 7). Just like we collect objects from region descriptions, we also collect the attributes attached to these objects. In Fig. 7, from the phrase “yellow fire hydrant,” we extract the attribute yellow for the fire hydrant. As with objects, we canonicalize all attributes to WordNet (Miller 1995); for example, yellow is mapped to `yellow.s.01` (of the color intermediate between green and orange in





**Fig. 8** Our dataset also captures the relationships and interactions between objects in our images. In this example, we show the relationship `jumping over` between the objects `man` and `fire hydrant`

the color spectrum; of something resembling the color of an egg yolk).

## 2.4 A Set of Relationships

Relationships connect two objects together. These relationships can be actions (e.g. `jumping over`), spatial (e.g. `is behind`), descriptive verbs (e.g. `wear`), prepositions (e.g. `with`), comparative (e.g. `taller than`), or prepositional phrases (e.g. `drive on`). For example, from the region description “man jumping over fire hydrant,” we extract the relationship `jumping over` between the objects `man` and `fire hydrant` (Fig. 8). These relationships are directed from one object, called the subject, to another, called the object. In this case, the subject is the man, who is performing the relationship `jumping over` on the object `fire hydrant`. Each relationship is canonicalized to a WordNet (Miller 1995) synset ID; i.e. `jumping` is canonicalized to `jump.a.1` (move forward by leaps and bounds). On average, each image in our dataset contains 21 relationships.

## 2.5 A Set of Region Graphs

Combining the objects, attributes, and relationships extracted from region descriptions, we create a directed graph representation for each of the regions. Examples of region graphs are shown in Fig. 4. Each region graph is a structured representation of a part of the image. The nodes in the graph represent objects, attributes, and relationships. Objects are linked to their respective attributes while relationships link one object to another. The links connecting two objects in

Fig. 4 point from the subject to the relationship and from the relationship to the other object.

## 2.6 One Scene Graph

While region graphs are localized representations of an image, we also combine them into a single scene graph representing the entire image (Fig. 3). The scene graph is the *union* of all region graphs and contains all objects, attributes, and relationships from each region description. By doing so, we are able to combine multiple levels of scene information in a more coherent way. For example in Fig. 4, the leftmost region description tells us that the “fire hydrant is yellow,” while the middle region description tells us that the “man is jumping over the fire hydrant.” Together, the two descriptions tell us that the “man is jumping over a yellow fire hydrant.”

## 2.7 A Set of Question Answer Pairs

We have two types of QA pairs associated with each image in our dataset: *freeform QAs*, based on the entire image, and *region-based QAs*, based on selected regions of the image. We collect 6 different types of questions per image: what, where, how, when, who, and why. In Fig. 4, “Q. What is the woman standing next to?; A. Her belongings” is a freeform QA. Each image has at least one question of each type listed above. Region-based QAs are collected by prompting workers with region descriptions. For example, we use the region “yellow fire hydrant” to collect the region-based QA: “Q. What color is the fire hydrant?; A. Yellow.” Region based QAs are based on the description and allow us to independently study how well models perform at answering questions using the image or the region description as input.

## 3 Related Work

We discuss existing datasets that have been released and used by the vision community for classification and object detection. We also mention work that has improved object and attribute detection models. Then, we explore existing work that has utilized representations similar to our relationships between objects. In addition, we dive into literature related to cognitive tasks like image description, question answering, and knowledge representation.

### 3.1 Datasets

Datasets (Table 1) have been growing in size as researchers have begun tackling increasingly complicated problems. *Caltech 101* (Fei-Fei et al. 2007) was one of the first datasets hand-curated for image classification, with 101 object cate-

**Table 1** A comparison of existing datasets with Visual Genome

	Images	Descriptions per image	Total objects	# Object categories	Objects per image	# Attributes categories	Attributes per image	# Relationship categories	Relationships per image	Question answers
YFCC100M (Thomee et al. 2016)	100,000,000	–	–	–	–	–	–	–	–	–
Tiny images (Torralba et al. 2008)	80,000,000	–	–	53,464	1	–	–	–	–	–
ImageNet (Deng et al. 2009)	14,197,122	–	14,197,122	21,841	1	–	–	–	–	–
ILSVRC detection (2012) (Russakovsky et al. 2015)	476,688	–	534,309	200	2.5	–	–	–	–	–
MS-COCO (Lin et al. 2014)	328,000	5	27,472	80	–	–	–	–	–	–
Flickr 30K (Young et al. 2014)	31,783	5	–	–	–	–	–	–	–	–
Caltech 101 (Fei-Fei et al. 2007)	9144	–	9144	102	1	–	–	–	–	–
Caltech 256 (Griffin et al. 2007)	30,608	–	30,608	257	1	–	–	–	–	–
Caltech pedestrian (Dollar et al. 2012)	250,000	–	350,000	1	1.4	–	–	–	–	–
Pascal detection (Everingham et al. 2010)	11,530	–	27,450	20	2.38	–	–	–	–	–
Abstract scenes (Zitnick and Parikh 2013)	10,020	–	58	11	5	–	–	–	–	–
aPascal (Farhadi et al. 2009)	12,000	–	–	–	–	64	–	–	–	–
Animal attributes (Lampert et al. 2009)	30,000	–	–	–	–	1,280	–	–	–	–
SUN attributes (Patterson et al. 2014)	14,000	–	–	–	–	700	700	–	–	–
Caltech birds (Wah et al. 2011)	11,788	–	–	–	–	312	312	–	–	–
COCO actions (Ronchi and Perona 2015)	10,000	–	–	–	5.2	–	–	156	20.7	–
Visual phrases (Sadeghi and Farhadi 2011)	–	–	–	–	–	–	–	17	1	–
VisKE (Sadeghi et al. 2015)	–	–	–	–	–	–	–	6500	–	–
DAQUAR (Malinowski and Fritz 2014)	1,449	–	–	–	–	–	–	–	–	12,468
COCO QA (Ren et al. 2015a)	123,287	–	–	–	–	–	–	–	–	117,684
Baidu (Gao et al. 2015)	120,360	–	–	–	–	–	–	–	–	250,569
VQA (Antol et al. 2015)	204,721	–	–	–	–	–	–	–	–	614,163
Visual Genome	108,077	50	3,843,636	33,877	35	68,111	26	42,374	21	1,773,258

We show that Visual Genome has an order of magnitude more descriptions and question answers. It also has a more diverse set of object, attribute, and relationship classes. Additionally, Visual Genome contains a higher density of these annotations per image. The number of distinct categories in Visual Genome are calculated by lower-casing and stemming names of objects, attributes and relationships

gories and 15–30 examples per category. One of the biggest criticisms of Caltech 101 was the lack of variability in its examples. *Caltech 256* (Griffin et al. 2007) increased the number of categories to 256, while also addressing some of the shortcomings of Caltech 101. However, it still had only a handful of examples per category, and most of its images contained only a single object. *LabelMe* (Russell et al. 2008) introduced a dataset with multiple objects per category. They also provided a web interface that experts and novices could use to annotate additional images. This web interface enabled images to be labeled with polygons, helping create datasets for image segmentation. The *Lotus Hill dataset* (Yao et al. 2007) contains a hierarchical decomposition of objects (vehicles, man-made objects, animals, etc.) along with segmentations. Only a small part of this dataset is freely available. *SUN* (Xiao et al. 2010), just like *LabelMe* (Russell et al. 2008) and *Lotus Hill* (Yao et al. 2007), was curated for object detection. Pushing the size of datasets even further, *80 Million Tiny Images* (Torralba et al. 2008) created a significantly larger dataset than its predecessors. It contains tiny (i.e.  $32 \times 32$  pixels) images that were collected using WordNet (Miller 1995) synsets as queries. However, because the data in *80 Million Images* were not human-verified, they contain numerous errors. *YFCC100M* (Thomee et al. 2016) is another large database of 100 million images that is still largely unexplored. It contains human generated and machine generated tags.

*Pascal VOC* (Everingham et al. 2010) pushed research from classification to object detection with a dataset containing 20 semantic categories in 11,000 images. *ImageNet* (Deng et al. 2009) took WordNet synsets and crowd-sourced a large dataset of 14 million images. They started the ILSVRC (Russakovsky et al. 2015) challenge for a variety of computer vision tasks. Together, ILSVRC and PASCAL provide a test bench for object detection, image classification, object segmentation, person layout, and action classification. *MS-COCO* (Lin et al. 2014) recently released its dataset, with over 328,000 images with sentence descriptions and segmentations of 80 object categories. The previous largest dataset for image-based QA, *VQA* (Antol et al. 2015), contains 204,721 images annotated with three question answer pairs. They collected a dataset of 614,163 freeform questions with 6.1M ground truth answers (10 per question) and provided a baseline approach in answering questions using an image and a textual question as the input.

*Visual Genome* aims to bridge the gap between all these datasets, collecting not just annotations for a large number of objects but also scene graphs, region descriptions, and question answer pairs for image regions. Unlike previous datasets, which were collected for a single task like image classification, the *Visual Genome* dataset was collected to be a general-purpose representation of the visual world, without bias toward a particular task. Our images contain an average

of 35 objects, which is almost an order of magnitude more dense than any existing vision dataset. Similarly, we contain an average of 26 attributes and 21 relationships per image. We also have an order of magnitude more unique objects, attributes, and relationships than any other dataset. Finally, we have 1.7 million question answer pairs, also larger than any other dataset for visual question answering.

### 3.2 Image Descriptions

One of the core contributions of *Visual Genome* is its descriptions for multiple regions in an image. As such, we mention other image description datasets and models in this subsection. Most work related to describing images can be divided into two categories: retrieval of human-generated captions and generation of novel captions. Methods in the first category use similarity metrics between image features from predefined models to retrieve similar sentences (Ordonez et al. 2011; Hodosh et al. 2013). Other methods map both sentences and their images to a common vector space (Ordonez et al. 2011) or map them to a space of triples (Farhadi et al. 2010). Among those in the second category, a common theme has been to use recurrent neural networks to produce novel captions (Kiros et al. 2014; Mao et al. 2014; Karpathy and Fei-Fei 2015; Vinyals et al. 2015; Chen and Lawrence Zitnick 2015; Donahue et al. 2015; Fang et al. 2015). More recently, researchers have also used a visual attention model (Xu et al. 2015).

One drawback of these approaches is their attention to describing only the most salient aspect of the image. This problem is amplified by datasets like Flickr 30K (Young et al. 2014) and MS-COCO (Lin et al. 2014), whose sentence descriptions tend to focus, somewhat redundantly, on these salient parts. For example, “an elephant is seen wandering around on a sunny day,” “a large elephant in a tall grass field,” and “a very large elephant standing alone in some brush” are 3 descriptions from the MS-COCO dataset, and all of them focus on the salient elephant in the image and ignore the other regions in the image. Many real-world scenes are complex, with multiple objects and interactions that are best described using multiple descriptions (Karpathy and Fei-Fei 2015; Lebrete et al. 2015). Our dataset pushes toward a more complete understanding of an image by collecting a dataset in which we capture not just scene-level descriptions but also myriad of low-level descriptions, the “grammar” of the scene.

### 3.3 Objects

Object detection is a fundamental task in computer vision, with applications ranging from identification of faces in photo software to identification of other cars by self-driving cars on the road. It involves classifying an object into a dis-



tinct category and localizing the object in the image. Visual Genome uses objects as a core component on which each visual scene is built. Early datasets include the face detection (Huang et al. 2008) and pedestrian datasets (Dollar et al. 2012). The PASCAL VOC and ILSVRC's detection dataset pushed research in object detection. But the images in these datasets are iconic and do not capture the settings in which these objects usually co-occur. To remedy this problem, MS-COCO (Lin et al. 2014) annotated real-world scenes that capture object contexts. However, MS-COCO was unable to describe all the objects in its images, since they annotated only 80 object categories. In the real world, there are many more objects than the ones captured by existing datasets. Visual Genome aims at collecting annotations for all visual elements that occur in images, increasing the number of distinct categories to 33,877.

### 3.4 Attributes

The inclusion of attributes allows us to describe, compare, and more easily categorize objects. Even if we haven't seen an object before, attributes allow us to infer something about it; for example, "yellow and brown spotted with long neck" likely refers to a giraffe. Initial work in this area involved finding objects with similar features (Malisiewicz et al. 2008) using exemplar SVMs. Next, textures were used to study objects (Varma and Zisserman 2005), while other methods learned to predict colors (Ferrari and Zisserman 2007). Finally, the study of attributes was explicitly demonstrated to lead to improvements in object classification (Farhadi et al. 2009). Attributes were defined to be parts (e.g. "has legs"), shapes (e.g. "spherical"), or materials (e.g. "furry") and could be used to classify new categories of objects. Attributes have also played a large role in improving fine-grained recognition (Goering et al. 2014) on fine-grained attribute datasets like CUB-2011 (Wah et al. 2011). In Visual Genome, we use a generalized formulation (Johnson et al. 2015), but we extend it such that attributes are not image-specific binaries but rather object-specific for each object in a real-world scene. We also extend the types of attributes to include size (e.g. "small"), pose (e.g. "bent"), state (e.g. "transparent"), emotion (e.g. "happy"), and many more.

### 3.5 Relationships

Relationship extraction has been a traditional problem in information extraction and in natural language processing. Syntactic features (Zhou et al. 2007; GuoDong et al. 2005), dependency tree methods (Culotta and Sorensen 2004; Bunescu and Mooney 2005), and deep neural networks (Socher et al. 2012; Zeng et al. 2014) have been employed to extract relationships between two entities in a sentence. However, in computer vision, very little work has

gone into learning or predicting relationships. Instead, relationships have been implicitly used to improve other vision tasks. Relative layouts between objects have improved scene categorization (Izadinia et al. 2014), and 3D spatial geometry between objects has helped object detection (Choi et al. 2013). Comparative adjectives and prepositions between pairs of objects have been used to model visual relationships and improved object localization (Gupta and Davis 2008).

Relationships have already shown their utility in improving visual cognitive tasks (Antol et al. 2014; Yang et al. 2012). A meaning space of relationships has improved the mapping of images to sentences (Farhadi et al. 2010). Relationships in a structured representation with objects have been defined as a graph structure called a *scene graph*, where the nodes are objects with attributes and edges are relationships between objects. This representation can be used to generate indoor images from sentences and also to improve image search (Chang et al. 2014; Johnson et al. 2015). We use a similar scene graph representation of an image that generalizes across all these previous works (Johnson et al. 2015). Recently, relationships have come into focus again in the form of question answering about associations between objects (Sadeghi et al. 2015). These questions ask if a relationship, involving generally two objects, is true, e.g. "do dogs eat ice cream?". We believe that relationships will be necessary for higher-level cognitive tasks (Johnson et al. 2015; Lu et al. 2016), so we collect the largest corpus of them in an attempt to improve tasks by actually understanding interactions between objects.

### 3.6 Question Answering

Visual question answering (QA) has been recently proposed as a proxy task of evaluating a computer vision system's ability to understand an image beyond object recognition and image captioning (Geman et al. 2015; Malinowski and Fritz 2014). Several visual QA benchmarks have been proposed in the last few months. The DAQUAR (Malinowski and Fritz 2014) dataset was the first toy-sized QA benchmark built upon indoor scene RGB-D images of NYU Depth v2 (Nathan Silberman and Fergus 2012). Most new datasets (Yu et al. 2015; Ren et al. 2015a; Antol et al. 2015; Gao et al. 2015) have collected QA pairs on MS-COCO images, either generated automatically by NLP tools (Ren et al. 2015a) or written by human workers (Yu et al. 2015; Antol et al. 2015; Gao et al. 2015).

In previous datasets, most questions concentrated on simple recognition-based questions about the salient objects, and answers were often extremely short. For instance, 90% of DAQUAR answers (Malinowski and Fritz 2014) and 89% of VQA answers (Antol et al. 2015) consist of single-word object names, attributes, and quantities. This limitation bounds their diversity and fails to capture the long-tail details

of the images. Given the availability of new datasets, an array of visual QA models have been proposed to tackle QA tasks. The proposed models range from SVM classifiers and probabilistic inference (Malinowski and Fritz 2014) to recurrent neural networks (Gao et al. 2015; Malinowski et al. 2015; Ren et al. 2015a) and convolutional networks (Ma et al. 2015). Visual Genome aims to capture the details of the images with diverse question types and long answers. These questions should cover a wide range of visual tasks from basic perception to complex reasoning. Our QA dataset of 1.7 million QAs is also larger than any currently existing dataset.

### 3.7 Knowledge Representation

A knowledge representation of the visual world is capable of tackling an array of vision tasks, from action recognition to general question answering. However, it is difficult to answer “what is the minimal viable set of knowledge needed to understand about the physical world?” (Hayes 1978). It was later proposed that there be a certain plurality to concepts and their related axioms (Hayes 1985). These efforts have grown to model physical processes (Forbus 1984) or to model a series of actions as scripts (Schank and Abelson 2013) for stories—both of which are not depicted in a single static image but which play roles in an image’s story (Vedantam et al. 2015b). More recently, NELL (Betteridge et al. 2009) learns probabilistic horn clauses by extracting information from the web. DeepQA (Ferrucci et al. 2010) proposes a probabilistic question answering architecture involving over 100 different techniques. Others have used Markov logic networks (Zhu et al. 2009; Niu et al. 2012) as their representation to perform statistical inference for knowledge base construction. Our work is most similar to that of those (Chen et al. 2013; Zhu et al. 2014, 2015; Sadeghi et al. 2015) who attempt to learn common-sense relationships from images. Visual Genome scene graphs can also be considered a *dense* knowledge representation for images. It is similar to the format used in knowledge bases in NLP.

## 4 Crowdsourcing Strategies

Visual Genome was collected and verified entirely by crowd workers from Amazon Mechanical Turk. In this section, we outline the pipeline employed in creating all the components of the dataset. Each component (region descriptions, objects, attributes, relationships, region graphs, scene graphs, questions and answers) involved multiple task stages. We mention the different strategies used to make our data accurate and to enforce diversity in each component. We also provide background information about the workers who helped make Visual Genome possible.

**Table 2** Geographic distribution of countries from where crowd workers contributed to Visual Genome

Country	Distribution (%)
United States	93.02
Philippines	1.29
Kenya	1.13
India	0.94
Russia	0.50
Canada	0.47
(Others)	2.65

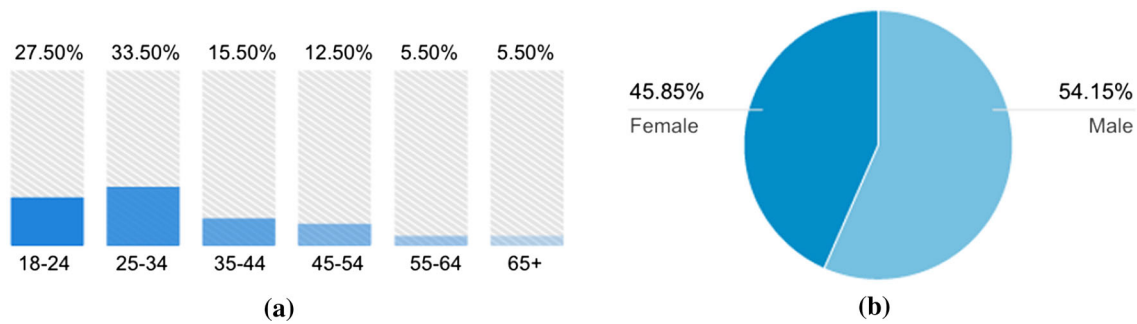
### 4.1 Crowd Workers

We used Amazon Mechanical Turk (AMT) as our primary source of annotations. Overall, a total of over 33, 000 unique workers contributed to the dataset. The dataset was collected over the course of 6 months after 15 months of experimentation and iteration on the data representation. Approximately 800, 000 Human Intelligence Tasks (HITs) were launched on AMT, where each HIT involved creating descriptions, questions and answers, or region graphs. Each HIT was designed such that workers manage to earn anywhere between \$6–\$8 per hour if they work continuously, in line with ethical research standards on Mechanical Turk (Salehi et al. 2015). Visual Genome HITs achieved a 94.1% retention rate, meaning that 94.1% of workers who completed one of our tasks went ahead to do more. Table 2 outlines the percentage distribution of the locations of the workers. 93.02% of workers contributed from the United States.

Figure 9a, b outline the demographic distribution of our crowd workers. This data was collected using a survey HIT. The majority of our workers were between the ages of 25 and 34 years old. Our youngest contributor was 18 years and the oldest was 68 years old. We also had a near-balanced split of 54.15% male and 45.85% female workers.

### 4.2 Region Descriptions

Visual Genome’s main goal is to enable the study of cognitive computer vision tasks. The next step towards understanding images requires studying relationships between objects in scene graph representations of images. However, we observed that collecting scene graphs directly from an image leads to workers annotating easy, frequently-occurring relationships like *wearing(man, shirt)* instead of focusing on salient parts of the image. This is evident from previous datasets (Johnson et al. 2015; Lu et al. 2016) that contain a large number of such relationships. After experimentation, we observed that when asked to describe an image using natural language, crowd workers naturally start with the most salient part of the image and then move to describing other



**Fig. 9** **a** Age and **b** gender distribution of Visual Genome's crowd workers

parts of the image one by one. Inspired by this finding, we focused our attention towards collecting a dataset of region descriptions that is diverse in content.

When a new image is added to the crowdsourcing pipeline with no annotations, it is sent to a worker who is asked to draw three bounding boxes and write three descriptions for the region enclosed by each box. Next, the image is sent to another worker along with the previously written descriptions. Workers are explicitly encouraged to write descriptions that have not been written before. This process is repeated until we have collected 50 region descriptions for each image. To prevent workers from having to skim through a long list of previously written descriptions, we only show them the top seven most similar descriptions. We calculate these most similar descriptions using BLEU-like (Papineni et al. 2002) (n-gram) scores between pairs of sentences. We define the similarity score  $S$  between a description  $d_i$  and a previous description  $d_j$  to be:

$$S_n(d_i, d_j) = b(d_i, d_j) \exp \left( \frac{1}{N} \sum_{n=1}^N \log p_n(d_i, d_j) \right) \quad (1)$$

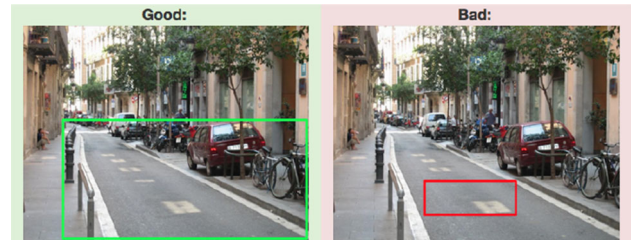
where we enforce a brevity penalty using:

$$b(d_i, d_j) = \begin{cases} 1 & \text{if } \text{len}(d_i) > \text{len}(d_j) \\ e^{1 - \frac{\text{len}(d_j)}{\text{len}(d_i)}} & \text{otherwise} \end{cases} \quad (2)$$

and  $p_n$  calculates the percentage of n-grams in  $d_i$  that match n-grams in  $d_j$ .

When a worker writes a new description, we programmatically enforce that it has not been repeated by using BLEU score thresholds set to 0.7 to ensure that it is dissimilar to descriptions from both of the following two lists:

1. **Image-Specific Descriptions** A list of all previously written descriptions for that image.
2. **Global Image Descriptions** A list of the top 100 most common written descriptions of all images in the dataset. This prevents very common phrases like “sky is blue”



**Fig. 10** Good (left) and bad (right) bounding boxes for the phrase “a street with a red car parked on the side,” judged on **coverage**

from dominating the set of region descriptions. The list of top 100 global descriptions is continuously updated as more data comes in.

Finally, we ask workers to draw bounding boxes that satisfy one requirement: **coverage**. The bounding box must cover all objects mentioned in the description. Figure 10 shows an example of a good box that covers both the street as well as the car mentioned in the description, as well as an example of a bad box.

### 4.3 Objects

Once 50 region descriptions are collected for an image, we extract the visual objects from each description. Each description is sent to one crowd worker, who extracts all the objects from the description and grounds each object as a bounding box in the image. For example, from Fig. 4, let's consider the description “woman in shorts is standing behind the man.” A worker would extract three objects: woman, shorts, and man. They would then draw a box around each of the objects. We require each bounding box to be drawn to satisfy two requirements: **coverage** and **quality**. Coverage has the same definition as described above in Sect. 4.2, where we ask workers to make sure that the bounding box covers the object completely (Fig. 11). Quality requires that each bounding box be as tight as possible around its object such that if the box's length or height were decreased by one pixel, it would no longer satisfy the coverage requirement.





**Fig. 11** Good (left) and bad (right) bounding boxes for the object `fox`, judged on both **coverage** as well as **quality**

Since a one pixel error can be physically impossible for most workers, we relax the definition of quality to four pixels.

Multiple descriptions for an image might refer to the same object, sometimes with different words. For example, a man in one description might be referred to as `person` in another description. We can thus use this crowdsourcing stage to build these co-reference chains. With each region description given to a worker to process, we include a list of previously extracted objects as suggestions. This allows a worker to choose a previously drawn box annotated as `man` instead of redrawing a new box for `person`.

Finally, to increase the speed with which workers complete this task, we also use Stanford’s dependency parser (Manning et al. 2014) to extract nouns automatically and send them to the workers as suggestions. While the parser manages to find most of the nouns, it sometimes misses compound nouns, so we avoided completely depending on this automated method. By combining the parser with crowdsourcing tasks, we were able to speed up our object extraction process without losing accuracy.

#### 4.4 Attributes, Relationships, and Region Graphs

Once all objects have been extracted from each region description, we can extract the attributes and relationships described in the region. We present each worker with a region description along with its extracted objects and ask them to add attributes to objects or to connect pairs of objects with relationships, based on the text of the description. From the description “woman in shorts is standing behind the man”, workers will extract the attribute `standing` for the woman and the relationships `in(woman, shorts)` and `behind(woman, man)`. Together, objects, attributes, and relationships form the region graph for a region description. Some descriptions like “it is a sunny day” do not contain any objects and therefore have no region graphs associated with them. Workers are asked to not generate any graphs for such descriptions. We create scene graphs by combining all the region graphs for an image by combining all the co-referenced objects from different region graphs.



**Fig. 12** Each object (`fox`) has only one bounding box referring to it (left). Multiple boxes drawn for the same object (right) are combined together if they have a minimum threshold of 0.9 intersection over union

#### 4.5 Scene Graphs

The scene graph is the union of all region graphs extracted from region descriptions. We merge nodes from region graphs that correspond to the same object; for example, `man` and `person` in two different region graphs might refer to the same object in the image. We say that objects from different graphs refer to the same object if their bounding boxes have an intersection over union of 0.9. However, this heuristic might contain false positives. So, before merging two objects, we ask workers to confirm that a pair of objects with significant overlap are indeed the same object. For example, in Fig. 12 (right), the `fox` might be extracted from two different region descriptions. These boxes are then combined together (Fig. 12, left) when constructing the scene graph.

#### 4.6 Questions and Answers

To create question answer (QA) pairs, we ask the AMT workers to write pairs of questions and answers about an image. To ensure quality, we instruct the workers to follow three rules: 1) start the questions with one of the “six Ws” (who, what, where, when, why and how); 2) avoid ambiguous and speculative questions; 3) be precise and unique, and relate the question to the image such that it is clearly answerable if and only if the image is shown.

We collected two separate types of QAs: freeform QAs and region-based QAs. In freeform QA, we ask a worker to look at an image and write eight QA pairs about it. To encourage diversity, we enforce that workers write at least three different Ws out of the six in their eight pairs. In region-based QA, we ask the workers to write a pair based on a given region. We select the regions that have large areas (more than 5k pixels) and long phrases (more than 4 words). This enables us to collect around twenty region-based pairs at the same cost of the eight freeform QAs. In general, freeform QA tends to yield more diverse QA pairs that enrich the question distribution; region-based QA tends to produce more factual QA pairs at a lower cost.

#### 4.7 Verification

All Visual Genome data go through a verification stage as soon as they are annotated. This stage helps eliminate incorrectly labeled objects, attributes, and relationships. It also helps remove region descriptions and questions and answers that might be correct but are vague (“This person seems to enjoy the sun.”), subjective (“room looks dirty”), or opinionated (“Being exposed to hot sun like this may cause cancer”).

Verification is conducted using two separate strategies: majority voting (Snow et al. 2008) and rapid judgments (Krishna et al. 2016). All components of the dataset except objects are verified using majority voting. Majority voting (Snow et al. 2008) involves three unique workers looking at each annotation and voting on whether it is factually correct. An annotation is added to our dataset if at least two (a majority) out of the three workers verify that it is correct.

We only use rapid judgments to speed up the verification of the objects in our dataset. Rapid judgments (Krishna et al. 2016) use an interface inspired by rapid serial visual processing that enable verification of objects with an order of magnitude increase in speed than majority voting.

#### 4.8 Canonicalization

All the descriptions and QAs that we collect are freeform worker-generated texts. They are not constrained by any limitations. For example, we do not force workers to refer to a man in the image as a man. We allow them to choose to refer to the man as person, boy, man, etc. This ambiguity makes it difficult to collect all instances of man from our dataset. In order to reduce the ambiguity in the concepts of our dataset and connect it to other resources used by the research community, we map all objects, attributes, relationships, and noun phrases in region descriptions and QAs to synsets in WordNet (Miller 1995). In the example above, person, boy, and man would map to the synsets: person.n.01 (a human being), male\_child.n.01 (a youthful male person) and man.n.03 (the generic use of the word to refer to any human being) respectively. Thanks to the WordNet hierarchy it is now possible to fuse those three expressions of the same concept into person.n.01 (a human being), which is the lowest common ancestor node of all aforementioned synsets.

We use the Stanford NLP tools (Manning et al. 2014) to extract the noun phrases from the region descriptions and QAs. Next, we map them to their most frequent matching synset in WordNet according to WordNet lexeme counts. We then refine this simple heuristic by hand-crafting mapping rules for the 30 most common failure cases. For example according to WordNet’s lexeme counts the most common semantic for “table” is table.n.01 (a set of data arranged in rows and columns). However in our

data it is more likely to see pieces of furniture and therefore bias the mapping towards table.n.02 (a piece of furniture having a smooth flat top that is usually supported by one or more vertical legs). The objects in our scene graphs are already noun phrases and are mapped to WordNet in the same way.

We normalize each attribute based on morphology (so called “stemming”) and map them to the WordNet adjectives. We include 15 hand-crafted rules to address common failure cases, which typically occur when the concrete or spatial sense of the word seen in an image is not the most common overall sense. For example, the synset long.a.02 (of relatively great or greater than average spatial extension) is less common in WordNet than long.a.01 (indicating a relatively great or greater than average duration of time), even though instances of the word “long” in our images are much more likely to refer to that spatial sense.

For relationships, we ignore all prepositions as they are not recognized by WordNet. Since the meanings of verbs are highly dependent upon their morphology and syntactic placement (e.g. passive cases, prepositional phrases), we try to find WordNet synsets whose sentence frames match with the context of the relationship. Sentence frames in WordNet are formalized syntactic frames in which a certain sense of a word might appear; e.g., play.v.01: participate in games or sport occurs in the sentence frames “Somebody [play]s” and “Somebody [play]s something.” For each verb-synset pair, we then consider the root hypernym of that synset to reduce potential noise from WordNet’s fine-grained sense distinctions. The WordNet hierarchy for verbs is segmented and originates from over 100 root verbs. For example, draw.v.01: cause to move by pulling traces back to the root hypernym move.v.02: cause to move or shift into a new position, while draw.v.02: get or derive traces to the root get.v.01: come into the possession of some thing concrete or abstract. We also include 20 hand-mapped rules, again to correct for WordNet’s lower representation of concrete or spatial senses.

These mappings are not perfect and still contain some ambiguity. Therefore, we send all our mappings along with the top four alternative synsets for each term to AMT. We ask workers to verify that our mapping was accurate and change the mapping to an alternative one if it was a better fit. We present workers with the concept we want to canonicalize along with our proposed corresponding synset with 4 additional options. To prevent workers from always defaulting to the our proposed synset, we do not explicitly specify which one of the 5 synsets presented is our proposed synset. Section 5.8 provides experimental precision and recall scores for our canonicalization strategy.

## 5 Dataset Statistics and Analysis

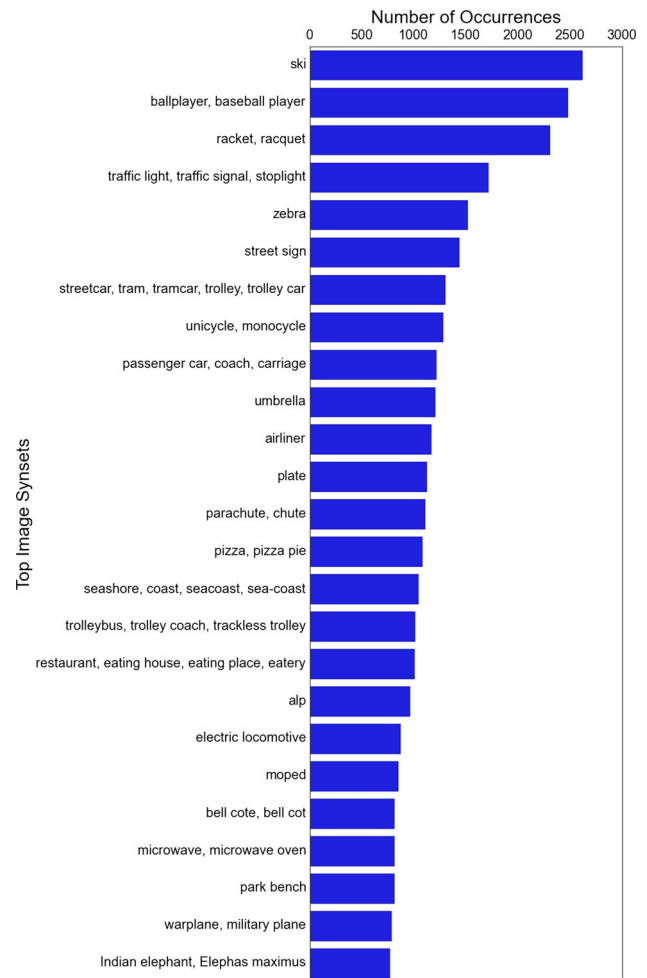
In this section, we provide statistical insights and analysis for each component of Visual Genome. Specifically, we examine the distribution of *images* (Sect. 5.1) and the collected data for *region descriptions* (Sect. 5.2) and *questions and answers* (Sect. 5.7). We analyze *region graphs* and *scene graphs* together in one section (Sect. 5.6), but we also break up these graph structures into their three constituent parts—*objects* (Sect. 5.3), *attributes* (Sect. 5.4), and *relationships* (Sect. 5.5)—and study each part individually. Finally, we describe our canonicalization pipeline and results (Sect. 5.8).

### 5.1 Image Selection

The Visual Genome dataset consists of all 108,077 creative commons images from the intersection of MS-COCO's (Lin et al. 2014) 328,000 images and YFCC100M's (Thomee et al. 2016) 100 million images. This allows Visual Genome annotations to be utilized together with the YFCC tags and MS-COCO's segmentations and full image captions. These images are real-world, non-iconic images that were uploaded onto Flickr by users. The images range from as small as 72 pixels wide to as large as 1280 pixels wide, with an average width of 500 pixels. We collected the WordNet synsets into which our 108,077 images can be categorized using the same method as ImageNet (Deng et al. 2009). Visual Genome images can be categorized into 972 ImageNet synsets. Note that objects, attributes and relationships are categorized separately into more than 18K WordNet synsets (Sect. 5.8). Figure 13 shows the top synsets to which our images belong. “ski” is the most common synset, with 2612 images; it is followed by “ballplayer” and “racket,” with all three synsets referring to images of people playing sports. Our dataset is somewhat biased towards images of people, as Fig. 13 shows; however, they are quite diverse overall, as the top 25 synsets each have over 800 images, while the top 50 synsets each have over 500 examples.

### 5.2 Region Description Statistics

One of the primary components of Visual Genome is its region descriptions. Every image includes an average of 50 regions with a bounding box and a descriptive phrase. Figure 14 shows an example image from our dataset with its 50 region descriptions. We display bounding boxes for only 6 out of the 50 descriptions in the figure to avoid clutter. These descriptions tend to be highly diverse and can focus on a single object, like in “A bag,” or on multiple objects, like in “Man taking a photo of the elephants.” They encompass the most salient parts of the image, as in “An elephant taking food from a woman,” while also capturing the background, as in “Small buildings surrounded by trees.”

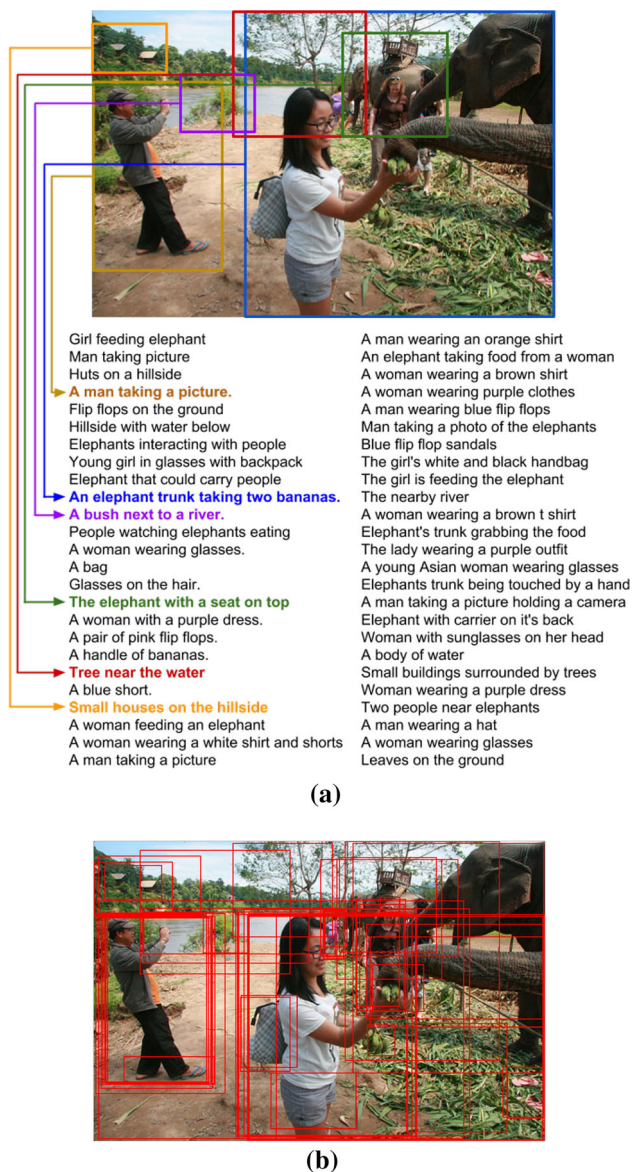


**Fig. 13** A distribution of the top 25 image synsets in the Visual Genome dataset. A variety of synsets are well represented in the dataset, with the top 25 synsets having at least 800 example images each. Note that an image synset is the label of the entire image according to the ImageNet ontology and are separate from the synsets for objects, attributes and relationships

MS-COCO (Lin et al. 2014) dataset is good at generating variations on a single scene-level descriptor. Consider three sentences from MS-COCO dataset on a similar image: “there is a person petting a very large elephant,” “a person touching an elephant in front of a wall,” and “a man in white shirt petting the cheek of an elephant.” These three sentences are single scene-level descriptions. In comparison, Visual Genome descriptions emphasize different regions in the image and thus are less semantically similar. To ensure diversity in the descriptions, we use BLEU score (Papineni et al. 2002) thresholds between new descriptions and all previously written descriptions. More information about crowdsourcing can be found in Sect. 4.

Region descriptions must be specific enough in an image to describe individual objects (e.g. “A bag”), but they must also be general enough to describe high-level concepts in an





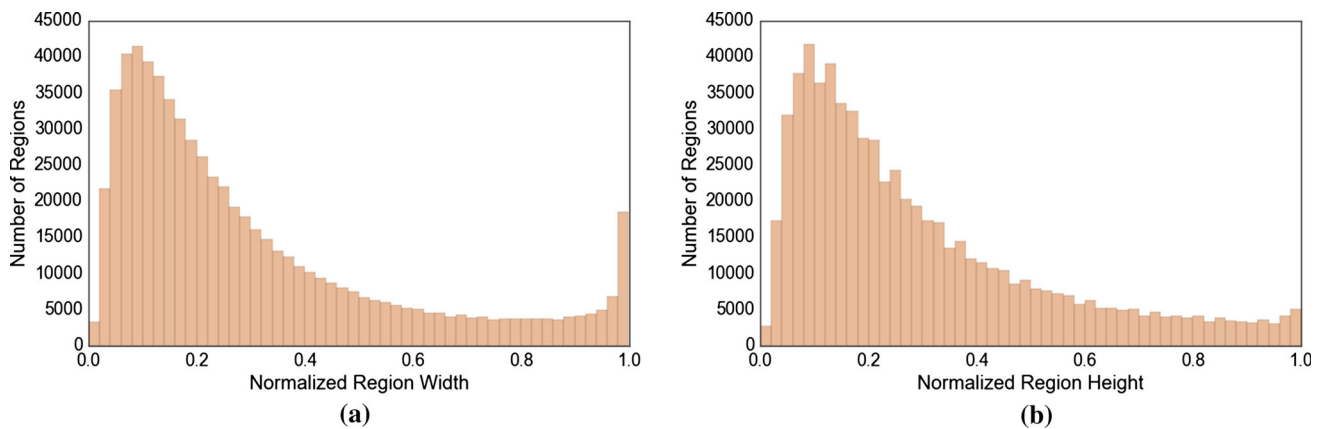
**Fig. 14** **a** An example image from the dataset with its region descriptions. We only display localizations for 6 of the 50 descriptions to avoid clutter; all 50 descriptions do have corresponding bounding boxes. **b** All 50 region bounding boxes visualized on the image

image (e.g. “A man being chased by a bear”). Qualitatively, we note that regions that cover large portions of the image tend to be general descriptions of an image, while regions that cover only a small fraction of the image tend to be more specific. In Fig. 15a, we show the distribution of regions over the width of the region normalized by the width of the image. We see that the majority of our regions tend to be around 10 to 15% of the image width. We also note that there are a large number of regions covering 100% of the image width. These regions usually include elements like “sky,” “ocean,” “snow,” “mountains,” etc. that cannot be bounded and thus span the entire image width. In Fig. 15b, we show a similar distribution

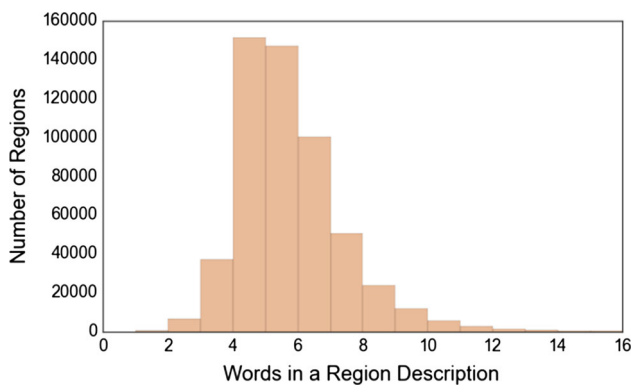
over the normalized height of the region. We see a similar overall pattern, as most of our regions tend to be very specific descriptions of about 10% to 15% of the image height. Unlike the distribution over width, however, we do not see a increase in the number of regions that span the entire height of the image, as there are no common visual equivalents that span images vertically. Out of all the descriptions gathered, only one or two of them tend to be global scene descriptions that are similar to MS-COCO (Lin et al. 2014) (Fig. 17).

In Fig. 16, we show the distribution of the length (word count) of these region descriptions. The average word count for a description is 5 words, with a minimum of 1 and a maximum of 12 words. In Fig. 18a, we plot the most common phrases occurring in our region descriptions, with common stop words removed. Common visual elements like “green grass,” “tree [in] distance,” and “blue sky” occur much more often than other, more nuanced elements like “fresh strawberry.” We also study descriptions with finer precision in Fig. 18b, where we plot the most common words used in descriptions. Again, we eliminate stop words from our study. Colors like “white” and “black” are the most frequently used words to describe visual concepts; we conduct a similar study on other captioning datasets including MS-COCO (Lin et al. 2014) and Flickr 30K (Young et al. 2014) and find a similar distribution with colors occurring most frequently. Besides colors, we also see frequent occurrences of common objects like “man” and “tree” and of universal visual elements like “sky.”

**Semantic Diversity** We also study the actual semantic contents of the descriptions. We use an unsupervised approach to analyze the semantics of these descriptions. Specifically, we use word2vec’s (Mikolov et al. 2013) pre-trained model on Google news corpus to convert each word in a description to a 300-dimensional vector. Next, we remove stop words and average the remaining words to get a vector representation of the whole region description. This pipeline is outlined in Fig. 17. We use hierarchical agglomerative clustering (Steinbach et al. 2000) on vector representations of each region description and find 71 semantic and syntactic groupings or “clusters.” Figure 19a shows four such example clusters. One cluster contains all descriptions related to tennis, like “A man swings the racquet” and “White lines on the ground of the tennis court,” while another cluster contains descriptions related to numbers, like “Three dogs on the street” and “Two people inside the tent.” To quantitatively measure the diversity of Visual Genome’s region descriptions, we calculate the number of clusters represented in a single image’s region descriptions. We show the distribution of the variety of descriptions for an image in Fig. 19b. We find that on average, each image contains descriptions from 17 different clusters. The image with the least diverse descriptions contains descriptions from 4 clusters, while the image



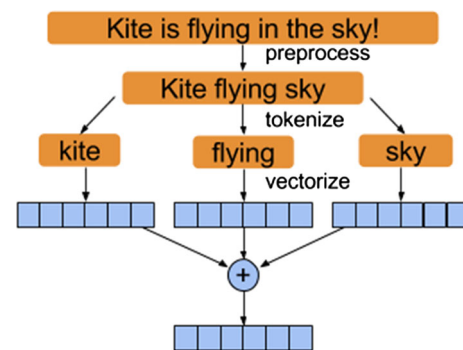
**Fig. 15** **a** A distribution of the width of the bounding box of a region description normalized by the image width. **b** A distribution of the height of the bounding box of a region description normalized by the image height



**Fig. 16** A distribution of the number of words in a region description. The average number of words in a region description is 5, with shortest descriptions of 1 word and longest descriptions of 16 words

with the most diverse descriptions contains descriptions from 26 clusters.

Finally, we also compare the descriptions in Visual Genome to the captions in MS-COCO. First we aggregate all Visual Genome and MS-COCO descriptions and remove all stop words. After removing stop words, the descriptions from both datasets are roughly the same length. We conduct a similar study, in which we vectorize the descriptions for each image and calculate each dataset's cluster diversity per image. We find that on average, 2 clusters are represented in the captions for each image in MS-COCO, with very few images in which 5 clusters are represented. Because each image in MS-COCO only contains 5 captions, it is not a fair comparison to compare the number of clusters represented in all the region descriptions in the Visual Genome dataset. We thus randomly sample 5 Visual Genome region descriptions per image and calculate the number of clusters in an image. We find that Visual Genome descriptions come from 4 or 5 clusters. We show our comparison results in Fig. 19c. The difference between the semantic diversity between the

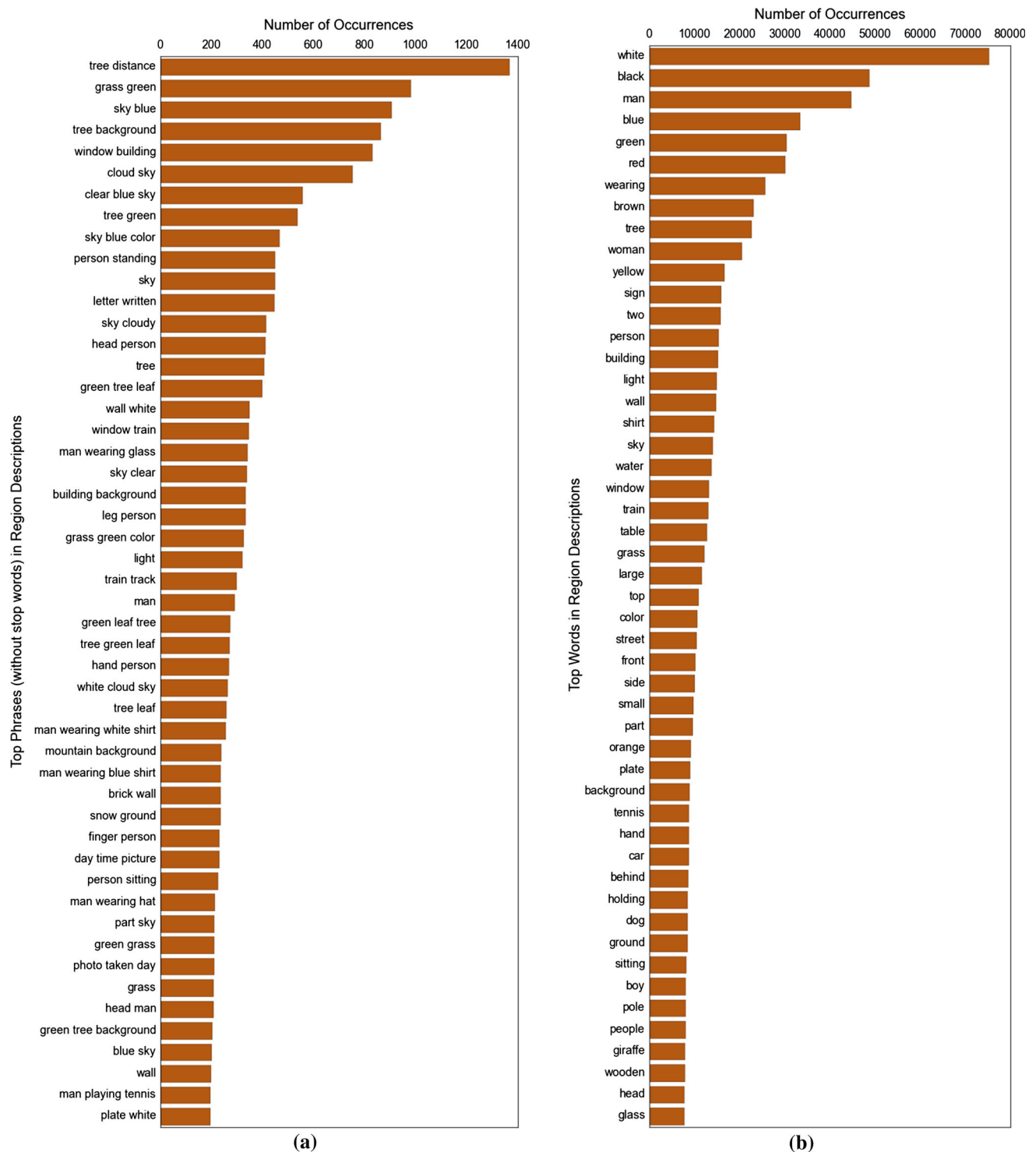


**Fig. 17** The process used to convert a region description into a 300-dimensional vectorized representation

two datasets is statistically significant ( $t = -240$ ,  $p < 0.01$ ) (Fig. 20).

### 5.3 Object Statistics

In comparison to related datasets, Visual Genome fares well in terms of object density and diversity (Table 3). Visual Genome contains approximately 35 objects per image, exceeding ImageNet (Deng et al. 2009), PASCAL (Everingham et al. 2010), MS-COCO (Lin et al. 2014), and other datasets by large margins. As shown in Fig. 21, there are more object categories represented in Visual Genome than in any other dataset. This comparison is especially pertinent with regards to Microsoft MS-COCO (Lin et al. 2014), which uses the same images as Visual Genome. The lower count of objects per category is a result of our higher number of categories. For a fairer comparison with ILSVRC 2014 Detection (Russakovsky et al. 2015), Visual Genome has about 2239 objects per category when only the top 200 categories are considered, which is comparable to ILSVRC's 2671.5 objects per category. For a fairer comparison with MS-COCO, Visual Genome has about 3768 objects per cat-

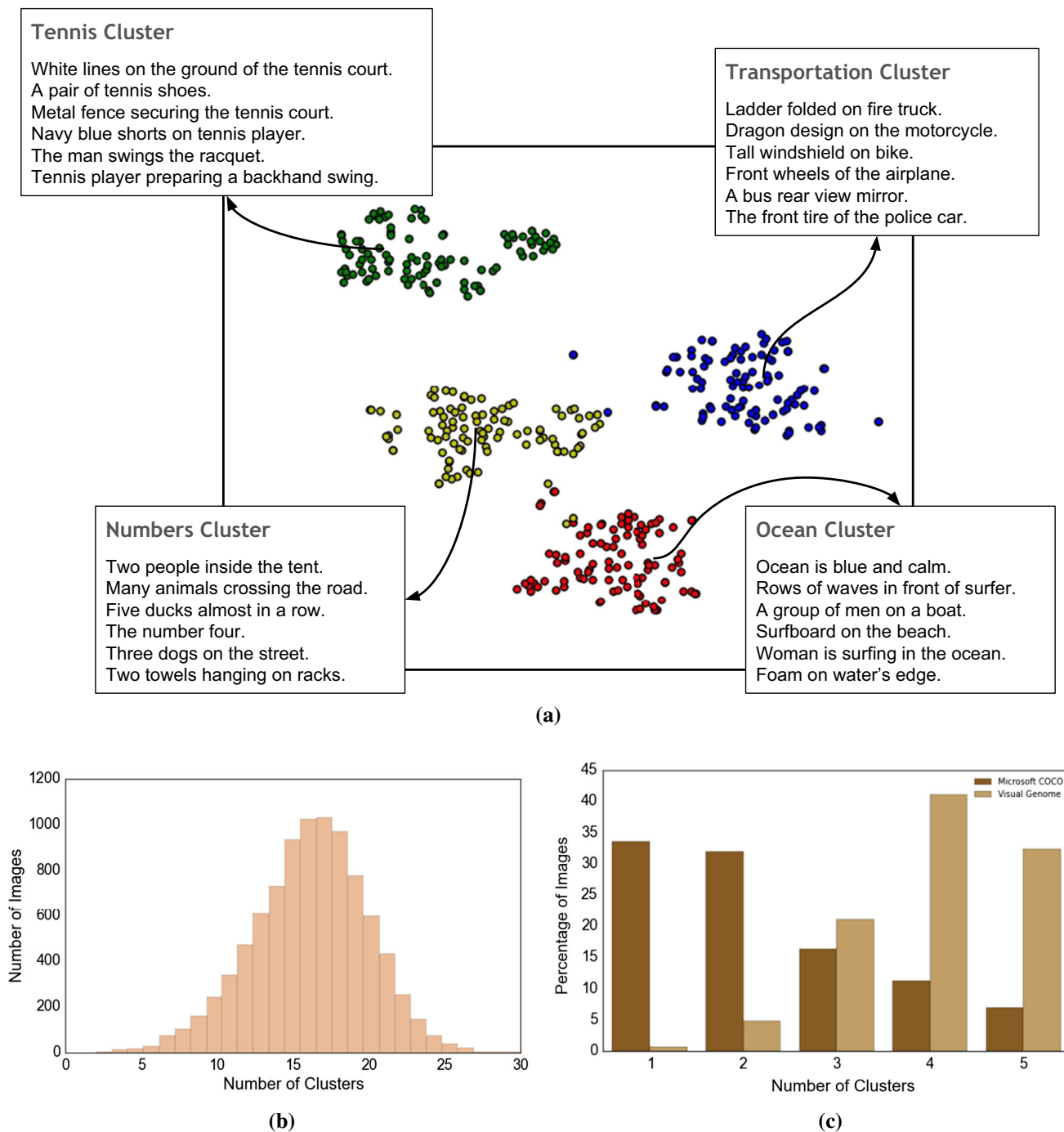


**Fig. 18** **a** A plot of the most common visual concepts or phrases that occur in region descriptions. The most common phrases refer to universal visual concepts like “blue sky,” “green grass,” etc. **b** A plot of the most frequently used words in region descriptions. Each word is treated

as an individual token regardless of which region description it came from. Colors occur the most frequently, followed by common objects like man and dog and universal visual concepts like “sky”

egory when only the top 80 categories are considered. This is comparable to MS-COCO’s (Lin et al. 2014) object distribution.

The 3,843,636 objects in Visual Genome come from a variety of categories. As shown in Fig. 22 (b), objects related to WordNet categories such as humans, animals, sports, and



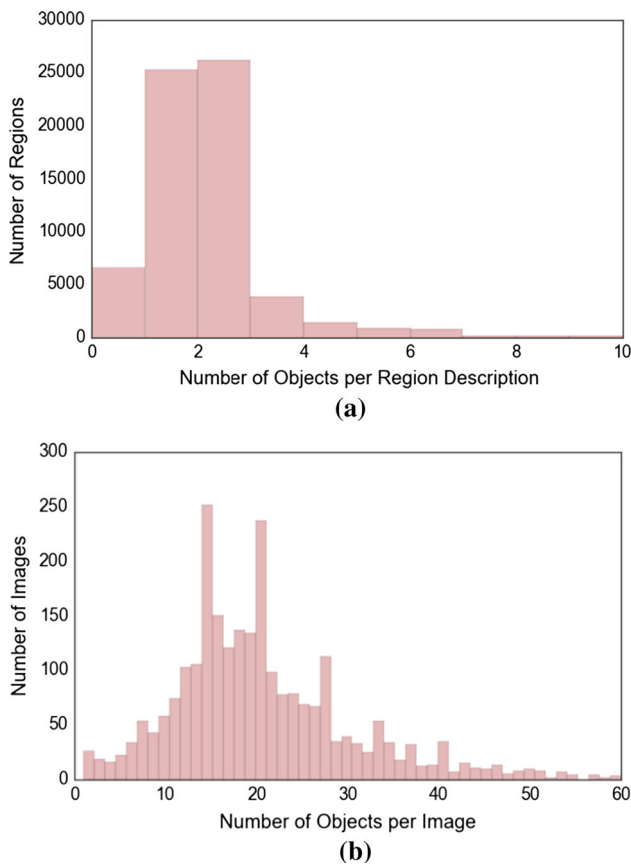
**Fig. 19** **a** Example illustration showing four clusters of region descriptions and their overall themes. Other clusters not shown due to limited space. **b** Distribution of images over number of clusters represented in each image's region descriptions. **c** We take Visual Genome with 5 random descriptions taken from each image and MS-COCO dataset

with all 5 sentence descriptions per image and compare how many clusters are represented in the descriptions. We show that Visual Genome's descriptions are more varied for a given image, with an average of 4 clusters per image, while MS-COCO's images have an average of 2 clusters per image

scenery are most common; this is consistent with the general bias in image subject matter in our dataset. Common objects like man, person, and woman occur especially frequently with occurrences of 24K, 17K, and 11K. Other objects that also occur in MS-COCO (Lin et al. 2014)

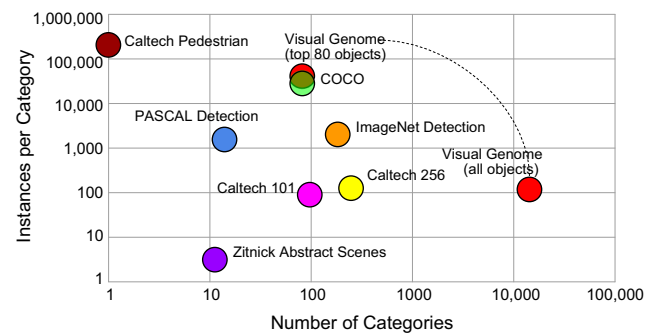
are also well represented with around 5000 instances on average. Figure 22a shows some examples of objects in images. Objects in Visual Genome span a diverse set of Wordnet categories like food, animals, and man-made structures.





**Fig. 20** **a** Distribution of the number of objects per region. Most regions have between 0 and 2 objects. **b** Distribution of the number of objects per image. Most images contain between 15 and 20 objects

It is important to look not only at what types of objects we have but also at the distribution of objects in images and regions. Figure 20a shows, as expected, that we have between 0 and 2 objects in each region on average. It is possible for regions to contain no objects if their descriptions refer to no explicit objects in the image. For example, a region described as “it is dark outside” has no objects to extract. Regions with only one object generally have descriptions that focus on the attributes of a single object. On the other hand, regions with two or more objects generally have descriptions that contain



**Fig. 21** Comparison of object diversity between various datasets. Visual Genome far surpasses other datasets in terms of number of categories. When considering only the top 80 object categories, it contains a comparable number of objects as MS-COCO. The *dashed line* is a visual aid connecting the two Visual Genome data points

both attributes of specific objects and relationships between pairs of objects.

As shown in Fig. 20b, each image contains on average around 35 distinct objects. Few images have an extremely high number of objects (e.g. over 40). Due to the image biases that exist in the dataset, we have twice as many annotations for men than we do of women.

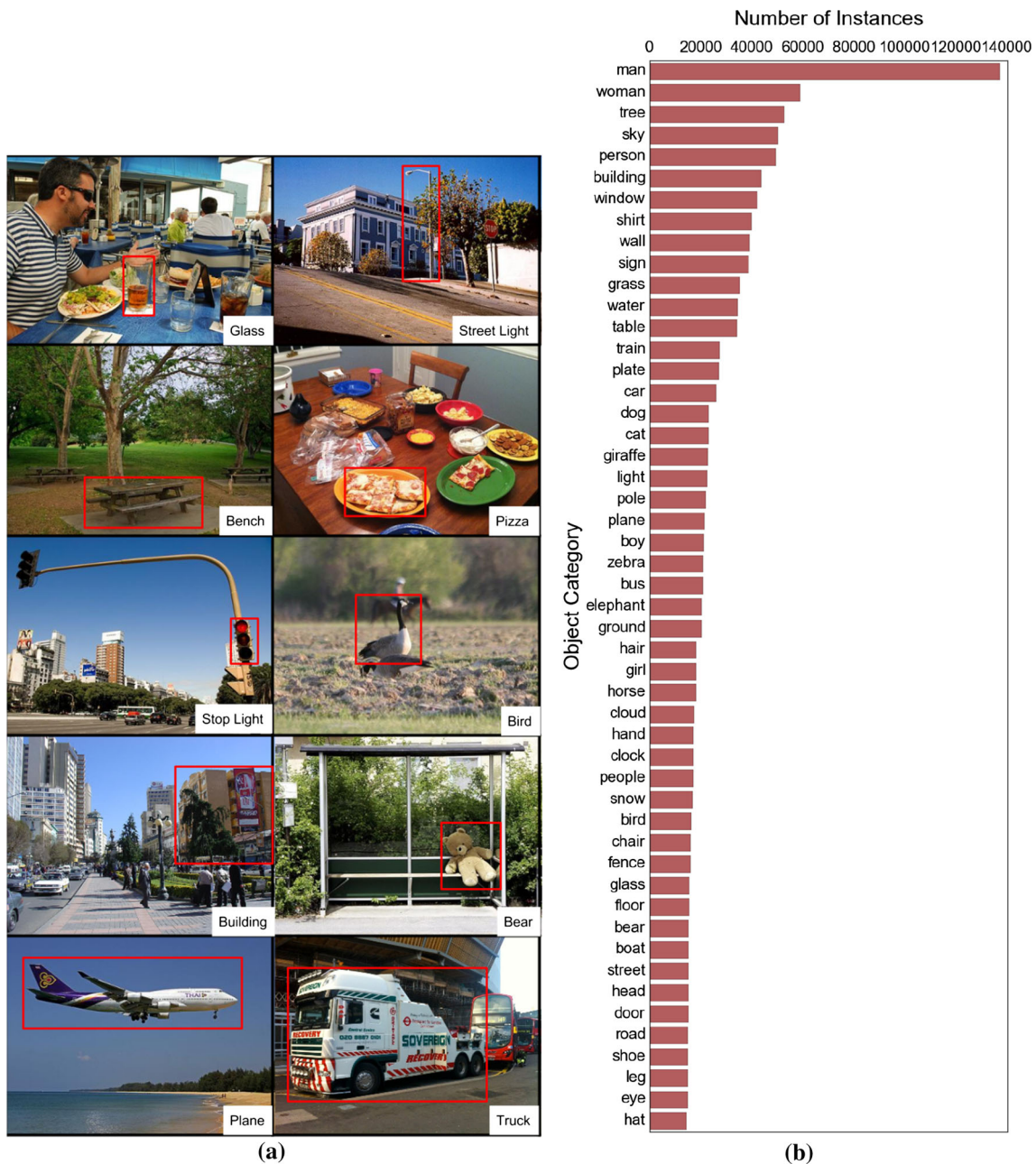
## 5.4 Attribute Statistics

Attributes allow for detailed description and disambiguation of objects in our dataset. Our dataset contains 2.8 million total attributes with 68,111 unique attributes. Attributes include colors (e.g. green), sizes (e.g. tall), continuous action verbs (e.g. standing), materials (e.g. plastic), etc. Each object can have multiple attributes.

On average, each image in Visual Genome contains 26 attributes (Fig. 23). Each region contains on average 1 attribute, though about 34% of regions contain no attribute at all; this is primarily because many regions are relationship-focused. Figure 24a shows the distribution of the most common attributes in our dataset. Colors (e.g. white, green) are by far the most frequent attributes. Also common are sizes (e.g. large) and materials (e.g. wooden). Figure 24b shows the distribution of attributes describing

**Table 3** Comparison of Visual Genome objects and categories to related datasets

	Visual Genome	ILSVRC det. (Russakovsky et al. 2015)	MS-COCO (Lin et al. 2014)	Caltech101 (Fei-Fei et al. 2007)	Caltech256 (Griffin et al. 2007)	PASCAL det. (Everingham et al. 2010)	Abstract scenes (Zitnick and Parikh 2013)
Images	108,077	476,688	328,000	9144	30,608	11,530	10,020
Total objects	3,843,636	534,309	2,500,000	9144	30,608	27,450	58
Total categories	33,877	200	80	102	257	20	11
Objects per category	113.45	2671.50	27472.50	90	119	1372.50	5.27



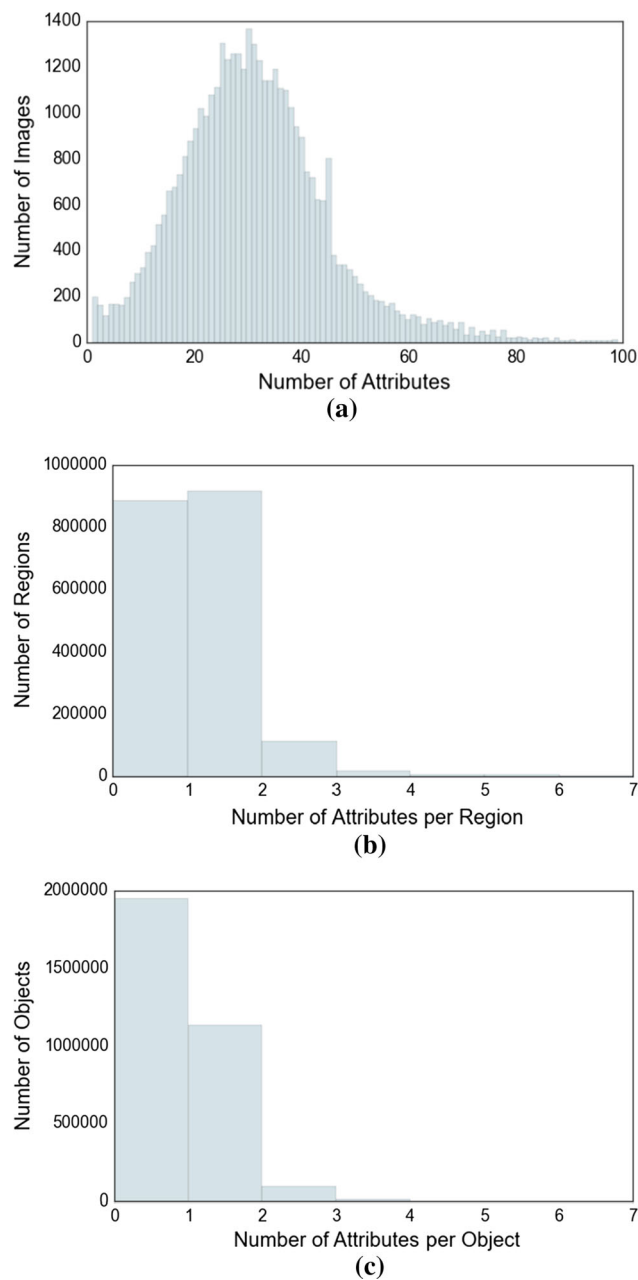
**Fig. 22** **a** Examples of objects in Visual Genome. Each object is localized in its image with a tightly drawn bounding box. **b** Plot of the most frequently occurring objects in images. People are the most frequently

occurring objects in our dataset, followed by common objects and visual elements like building, shirt, and sky

people (e.g. man, girls, and person). The most common attributes describing people are intransitive verbs describing their states of motion (e.g. standing and walking). Certain sports (e.g. skiing, surfboarding) are over-represented due to an image bias towards these sports.

**Attribute Graphs** We also qualitatively analyze the attributes in our dataset by constructing co-occurrence graphs, in which

nodes are unique attributes and edges connect those attributes that describe the same object. For example, if an image contained a “large black dog” (large(dog), black(dog)) and another image contained a “large yellow cat” (large(cat), yellow(cat)), its attributes would form an incomplete graph with edges (large, black) and (large, yellow). We create two such graphs: one for both the total set of attributes and a second where we consider only objects that



**Fig. 23** Distribution of the number of attributes **a** per image, **b** per region description, **c** per object

refer to people. A subgraph of the 16 most frequently connected (co-occurring) person-related attributes is shown in Fig. 25a.

Cliques in these graphs represent groups of attributes in which at least one co-occurrence exists for each pair of attributes. In the previous example, if a third image contained a “black and yellow taxi” (black(taxi), yellow(taxi)), the resulting third edge would create a clique between the attributes black, large, and yellow. When calculated across the entire Visual Genome dataset, these cliques pro-

vide insight into commonly perceived traits of different types of objects. Figure 25b is a selected representation of three example cliques and their overlaps. From just a clique of attributes, we can predict what types of objects are usually referenced. In Fig. 25b, we see that these cliques describe an animal (left), water body (top right), and human hair (bottom right).

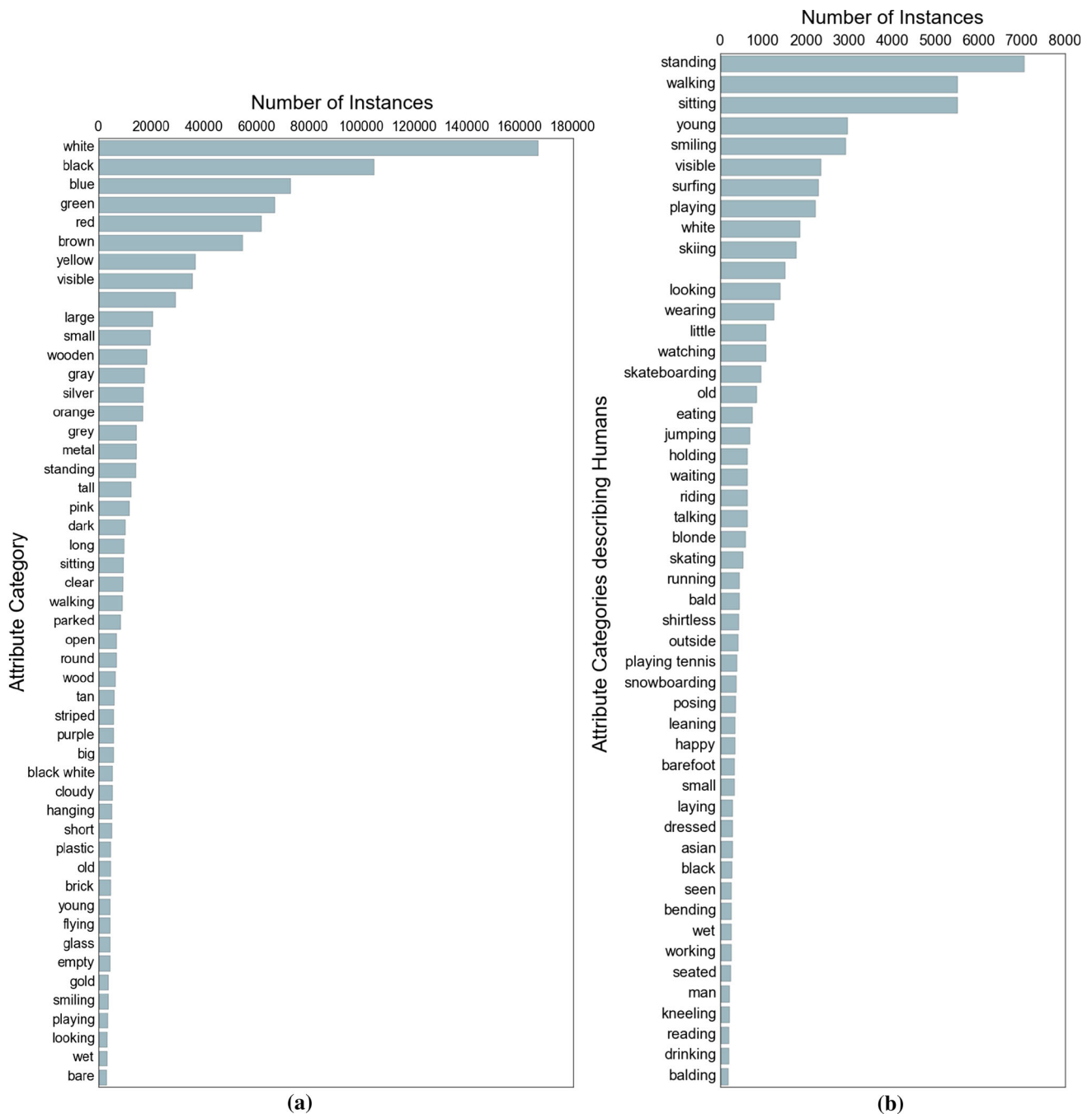
Other cliques (not shown) can also uniquely identify object categories. In our set, one clique contains athletic, young, fit, skateboarding, focused, teenager, male, skinny, and happy, capturing some of the common traits of skateboarders in our set. Another such clique has shiny, small, metal, silver, rusty, parked, and empty, most likely describing a subset of cars. From these cliques, we can thus infer distinct objects and object types based solely on their attributes, potentially allowing for highly specific object identification based on selected characteristics.

### 5.5 Relationship Statistics

Relationships are the core components that link objects in our scene graphs. Relationships are directional, i.e. they involve two objects, one acting as the subject and one as the object of a predicate relationship. We denote all relationships in the form *relationship(subject, object)*. For example, if a man is swinging a bat, we write *swinging(man, bat)*. Relationships can be spatial (e.g. *inside\_of*), action (e.g. *swinging*), compositional (e.g. *part\_of*), etc. More complex relationships such as *standing\_on*, which includes both an action and a spatial aspect, are also represented. Relationships are extracted from region descriptions by crowd workers, similarly to attributes and objects. Visual Genome contains a total of 42,374 unique relationships, with over 2,347,187 million total relationships.

Figure 26a shows the distribution of relationships per region description. On average, we have 1 relationship per region, with a maximum of 7. We also have some descriptions like “an old, tall man,” which have multiple attributes associated with the man but no relationships. Figure 26b is a distribution of relationships per image object. Finally, Fig. 26c shows the distribution of relationships per image. Each image has an average of 19 relationships, with a minimum of 1 relationship and with a maximum of over 80 relationships.

**Top Relationship Distributions** We display the most frequently occurring relationships in Fig. 27a. on is the most common relationship in our dataset. This is primarily because of the flexibility of the word on, which can refer to spatial configuration (on top of), attachment (hanging on), etc. Other common relationships involve actions like holding and wearing and spatial configu-



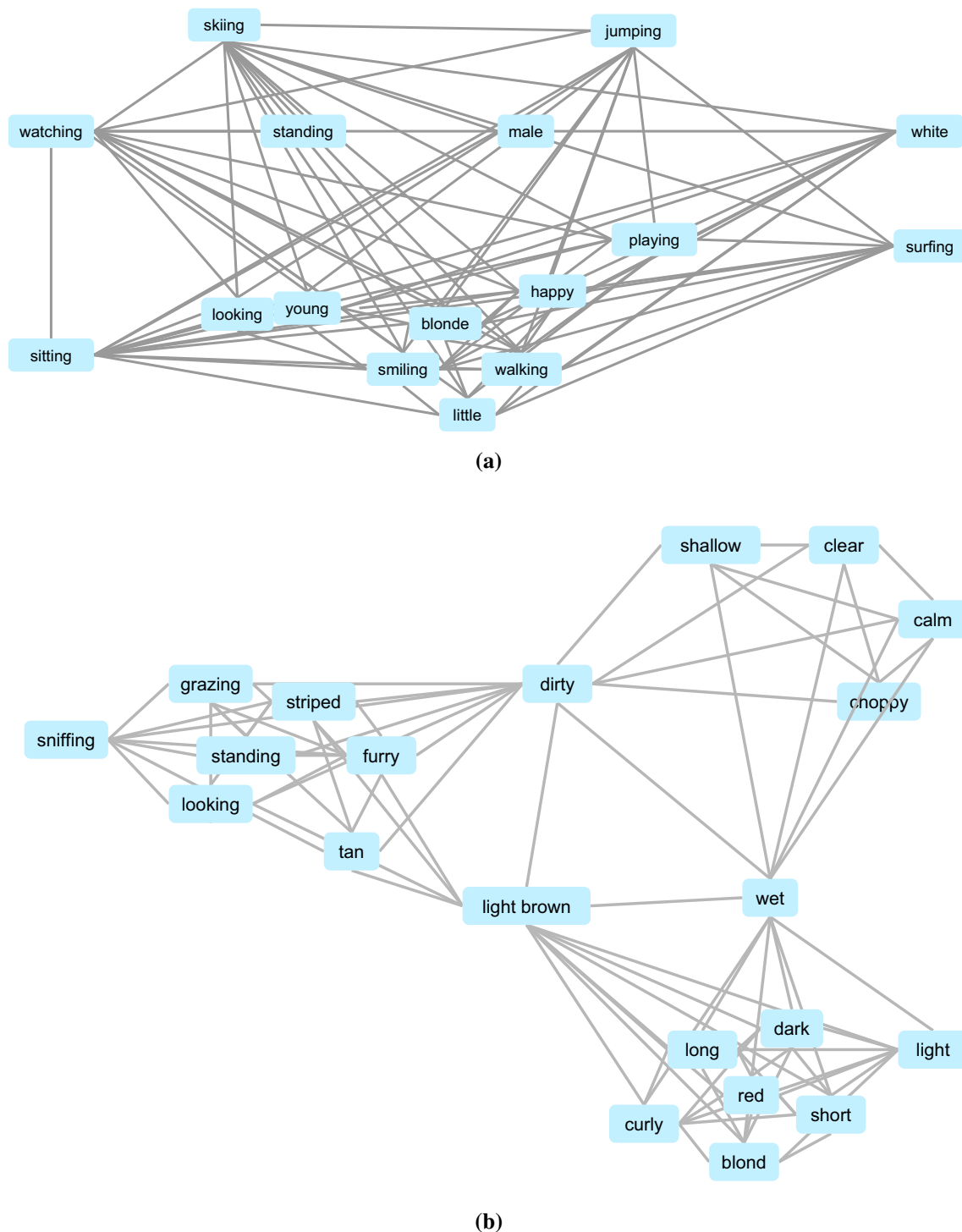
**Fig. 24** **a** Distribution showing the most common attributes in the dataset. Colors (e.g. white, red) and materials (e.g. wooden, metal) are the most common. **b** Distribution showing the number of attributes describing people. State-of-motion verbs (e.g. standing,

walking) are the most common, while certain sports (e.g. skiing, surfing) are also highly represented due to an image source bias in our image set

rations like behind, next to, and under. Figure 27b shows a similar distribution but for relationships involving people. Here we notice more human-centric relationships or actions such as kissing, chatting with, and talking to. The two distributions follow a Zipf distribution.

*Understanding Affordances Relationships* allow us to also understand the affordances of objects. Figure 28a shows the distribution for subjects while Fig. 28b shows a similar distribution for objects. Comparing the two, we find clear patterns of people-like subject entities such as person, man, policeman, boy, and skateboarder that can



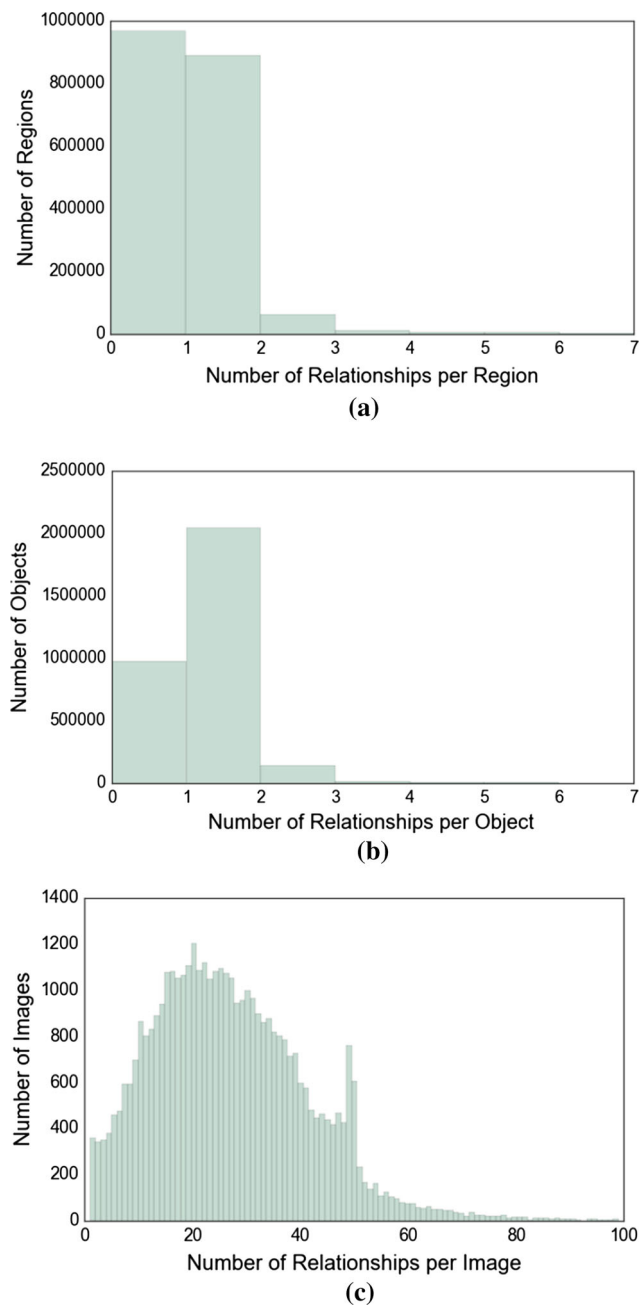


**Fig. 25** **a** Graph of the person-describing attributes with the most co-occurrences. *Edge thickness* represents the frequency of co-occurrence of the two nodes. **b** A subgraph showing the co-occurrences and inter-

sections of three cliques, which appear to describe water (*top right*), hair (*bottom right*), and some type of animal (*left*). Edges between cliques have been removed for clarity

ride other objects; the other distribution contains objects that afford riding, such as horse, bike, elephant, motorcycle, and skateboard. We can also learn spe-

cific common-sense knowledge, like that zebras eat hay and grass while a person eats pizzas and burgers and that couches usually have pillows on them.



**Fig. 26** Distribution of relationships **a** per image region, **b** per image object, **c** per image

**Related Work Comparison** It is also worth mentioning in this section some prior work on relationships. The concept of visual relationships has already been explored in Visual Phrases (Sadeghi and Farhadi 2011), who introduced a dataset of 17 such relationships such as *next\_to(person, bike)* and *riding(person, horse)*. However, their dataset is limited to just these 17 relationships. Similarly, the MS-COCO-a scene graph dataset (Ronchi and Perona 2015) introduced 156 actions that humans performed in MS-COCO’s dataset (Lin et al. 2014). They show that to exhaustively

describe “common” images involving humans, only a small set of visual actions is needed. However, their dataset is limited to just actions, while our relationships are more general and numerous, with over 42,374 unique relationships. Finally, VisKE (Sadeghi et al. 2015) introduced 6500 relationships, but in a much smaller dataset of images than Visual Genome.

## 5.6 Region and Scene Graph Statistics

We introduce in this paper the largest dataset of scene graphs to date. We use these graph representations of images as a deeper understanding of the visual world. In this section, we analyze the properties of these representations, both at the region-level through region graphs and at the image level through scene graphs. We also briefly explore other datasets with scene graphs and provide aggregate statistics on our entire dataset.

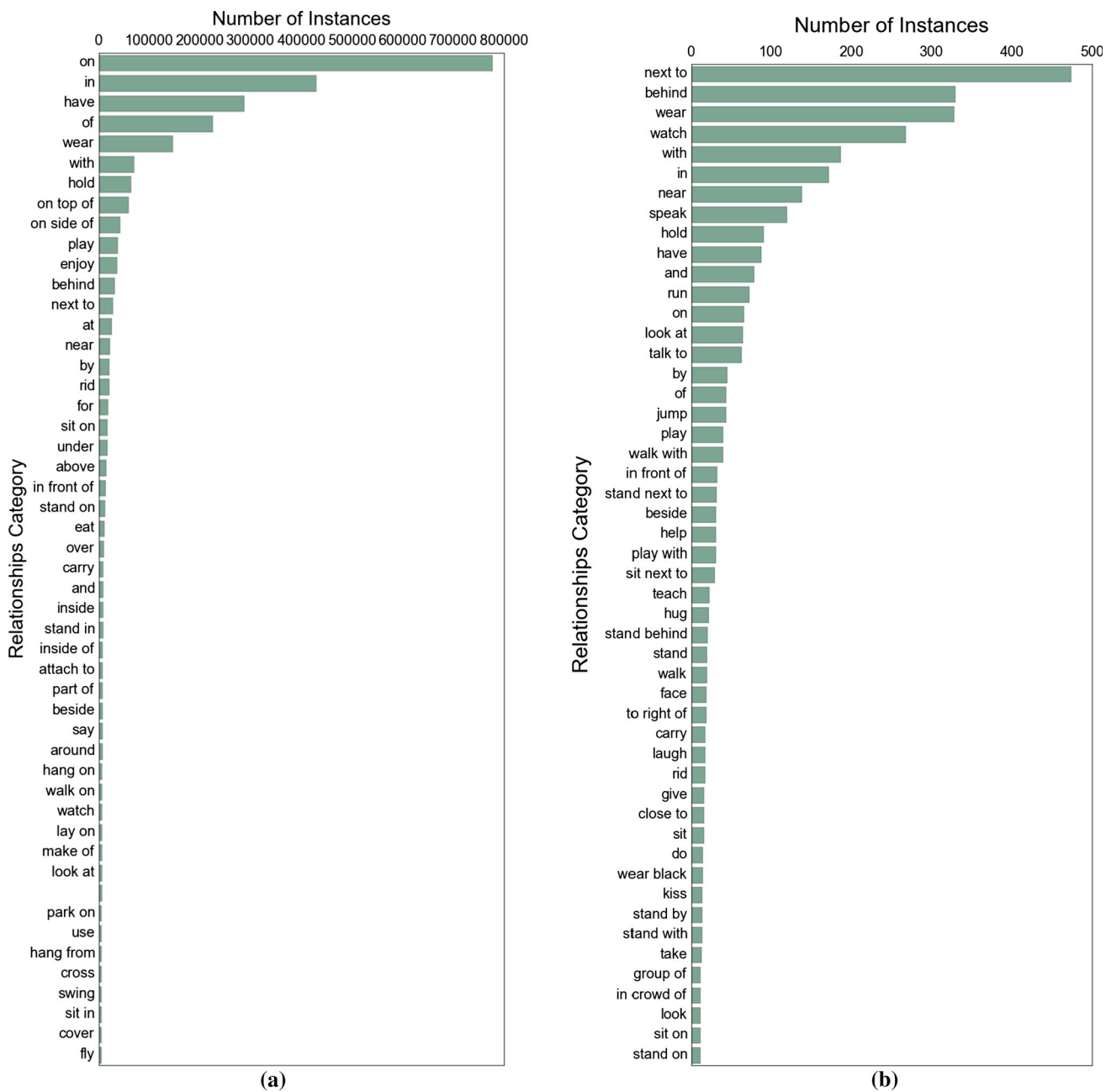
In previous work, scene graphs have been collected by asking humans to write a list of triples about an image (Johnson et al. 2015). However, unlike them, we collect graphs at a much more fine-grained level: the region graph. We obtained our graphs by asking workers to create them from the descriptions we collected from our regions. Therefore, we end up with multiple graphs for an image, one for every region description. Together, we can combine all the individual region graphs to aggregate a scene graph for an image. This scene graph is made up of all the individual region graphs. In our scene graph representation, we merge all the objects that referenced by multiple region graphs into one node in the scene graph.

Each of our images has between 5 to 100 region graphs per image, with an average of 50. Each image has exactly one scene graph. Note that the number of region descriptions and the number of region graphs for an image are not the same. For example, consider the description “it is a sunny day”. Such a description contains no objects, which are the building blocks of a region graph. Therefore, such descriptions have no region graphs associated with them.

Objects, attributes, and relationships occur as a normal distribution in our data. Table 4 shows that in a region graph, there are an average of 0.71 objects, 0.52 attributes, and 21 relationships. Each scene graph and consequently each image has average of 35 objects, 26 attributes, and 21 relationships.

## 5.7 Question Answering Statistics

We collected 1,773,258 question answering (QA) pairs on the Visual Genome images. Each pair consists of a question and its correct answer regarding the content of an image. On average, every image has 17 QA pairs. Rather than collecting unconstrained QA pairs as previous work has done (Antol et al. 2015; Gao et al. 2015; Malinowski and Fritz 2014), each

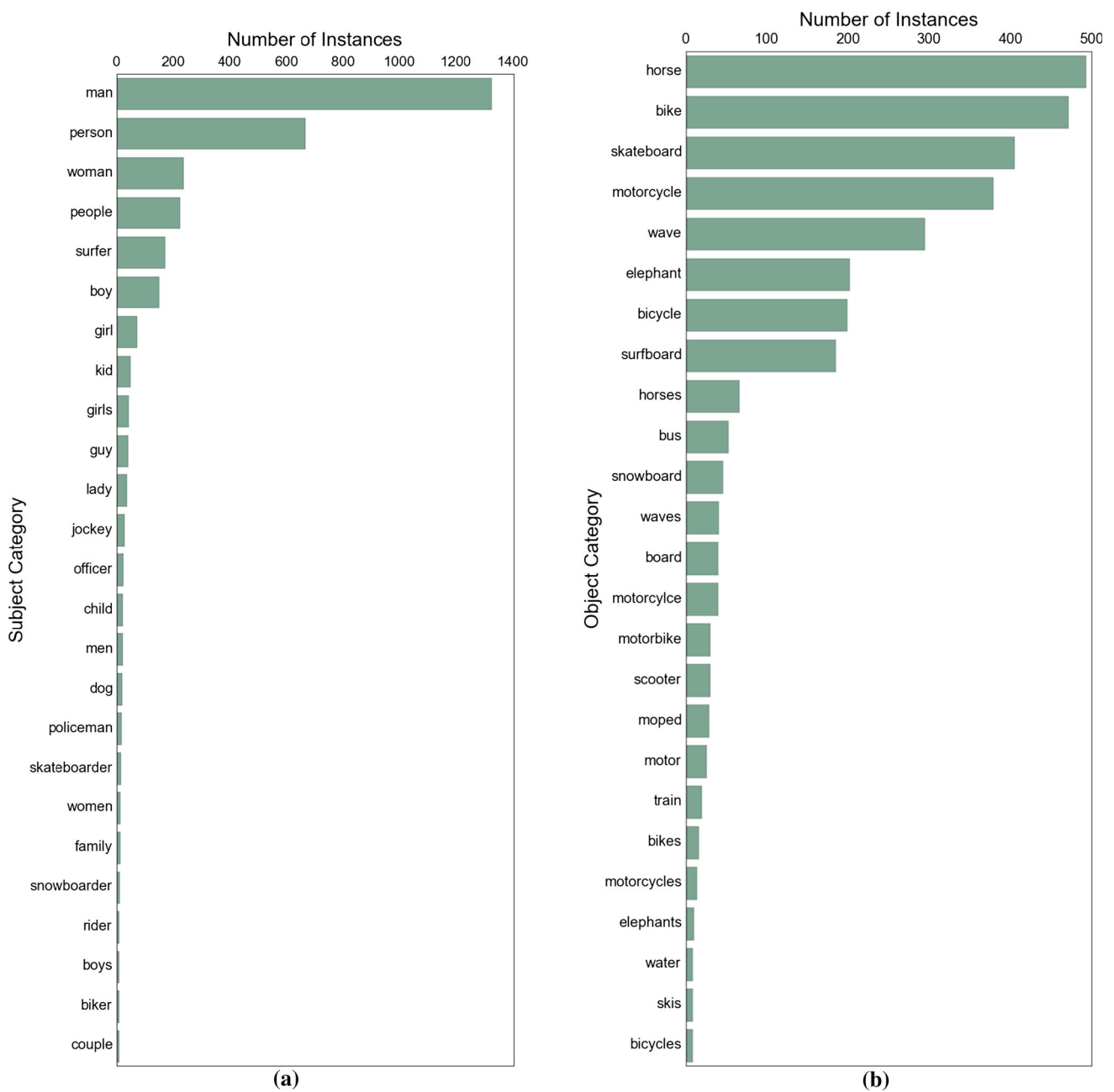


**Fig. 27** **a** A sample of the most frequent relationships in our dataset. In general, the most common relationships are spatial (on top of, on side of, etc.). **b** A sample of the most frequent relationships involv-

ing humans in our dataset. The relationships involving people tend to be more action oriented (walk, speak, run, etc.)

question in Visual Genome starts with one of the six Ws – what, where, when, who, why, and how. There are two major benefits to focusing on six types of questions. First, they offer a considerable coverage of question types, ranging from basic perceptual tasks (e.g. recognizing objects and scenes) to complex common sense reasoning (e.g. inferring motivations of people and causality of events). Second, these categories present a natural and consistent stratification of task diffi-

culty, indicated by the baseline performance in Sect. 6.4. For instance, *why* questions that involve complex reasoning lead to the poorest performance (3.4% top-100 accuracy compared to 9.6% top-100 accuracy of the next lowest) of the six categories. This enables us to obtain a better understanding of the strengths and weaknesses of today's computer vision models, which sheds light on future directions in which to proceed.



**Fig. 28** **a** Distribution of subjects for the relationship *riding*. **b** Distribution of objects for the relationship *riding*. Subjects comprise of people-like entities like *person*, *man*, *policeman*, *boy*, and

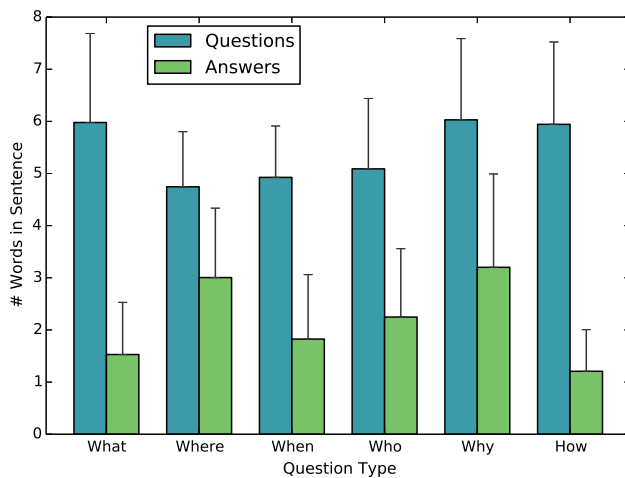
*skateboarder* that can ride other objects. On the other hand, objects like *horse*, *bike*, *elephant* and *motorcycle* are entities that can afford riding

We now analyze the diversity and quality of our questions and answers. Our goal is to construct a large-scale visual question answering dataset that covers a diverse range of question types, from basic cognition tasks to complex reasoning tasks. We demonstrate the richness and diversity of our QA pairs by examining the distributions of questions and answers in Fig. 29.

**Question Type Distributions** The questions naturally fall into the 6W categories via their interrogative words. Inside each of the categories, the second and following words categorize the questions with increasing granularity. Inspired by VQA (Antol et al. 2015), we show the distributions of the questions by their first three words in Fig. 30. We can see that “what” is the most common of the six categories. A notable





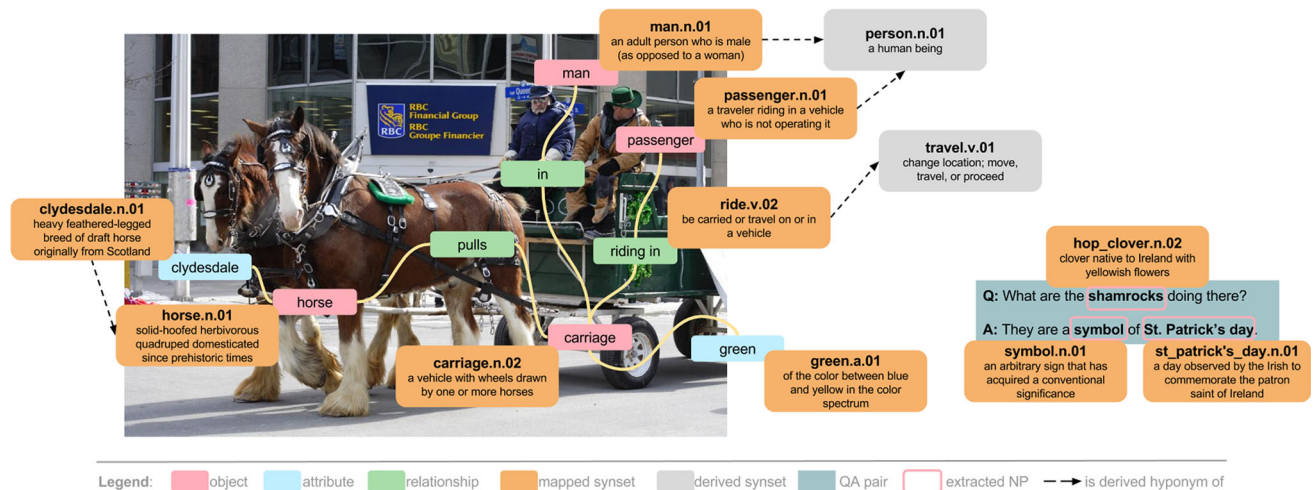


**Fig. 31** Question and answer lengths by question type. The bars show the average question and answer lengths of each question type. The whiskers show the standard deviations. The factual questions, such as “what” and “how” questions, usually come with short answers of a single object or a number. This is only because “how” questions are disproportionately counting questions that start with “how many”. Questions from the “where” and “why” categories usually have phrases and sentences as answers

relationships, and attributes in Visual Genome. By “canonicalization,” we refer to word sense disambiguation (WSD) by mapping the components in our dataset to their respective synsets in the WordNet ontology (Miller 1995). This mapping reduces the noise in the concepts contained in the dataset and also facilitates the linkage between Visual Genome and other data sources such as ImageNet (Deng et al. 2009), which is built on top of the WordNet ontology.

Figure 32 shows an example image from the Visual Genome dataset with its components canonicalized. For example, horse is canonicalized as horse.n.01: solid-hoofed herbivorous quadruped domesticated since prehistoric times. Its attribute, clydesdale, is canonicalized as its breed clydesdale.n.01: heavy feathered-legged breed of draft horse originally from Scotland. We also show an example of a QA from which we extract the nouns shamrocks, symbol, and St. Patrick’s day, all of which we canonicalize to WordNet as well.

**Related Work** Canonicalization, or WSD (Pal and Saha 2015), has been used in numerous applications, including machine translation, information retrieval, and information extraction (Rothe and Schütze 2015; Leacock et al. 1998). In English sentences, sentences like “He scored a goal” and “It was his goal in life” carry different meanings for the word “goal.” Understanding these differences is crucial for translating languages and for returning correct results for a query. Similarly, in Visual Genome, we ensure that all our components are canonicalized to understand how different objects are related to each other; for example, “person” is a hypernym of “man” and “woman.” Most past canonicalization models use precision, recall, and F1 score to evaluate on the Semeval dataset (Mihalcea et al. 2004). The current state-of-the-art performance on Semeval is an F1 score of 75.8% (Chen et al. 2014). Since our canonicalization setup is different from the Semeval benchmark (we have an open vocabulary and no annotated ground truth for evaluation), our canonicalization



**Fig. 32** An example image from the Visual Genome dataset with its region descriptions, QA pairs, objects, attributes, and relationships canonicalized. The large text boxes are WordNet synsets referenced by this image. For example, the carriage is mapped to carriage.n.02: a vehicle with wheels drawn by

one or more horses. We do not show the bounding boxes for the objects in order to allow readers to see the image clearly. We also only show a subset of the scene graph for this image to avoid cluttering the figure

**Table 5** Precision, recall, and mapping accuracy percentages for object, attribute, and relationship canonicalization

	Precision	Recall
Objects	88.0	98.5
Attributes	85.7	95.9
Relationships	92.9	88.5

method is not directly comparable to these existing methods. We do however, achieve a similar precision and recall score on a held-out test set described below (Table 5).

**Region Descriptions and QAs** We canonicalize all objects mentioned in all region descriptions and QA pairs. Because objects need to be extracted from the phrase text, we use Stanford NLP tools (Manning et al. 2014) to extract the noun phrases in each region description and QA, resulting in 99% recall of noun phrases from a subset of 200 region descriptions we manually annotated. After obtaining the noun phrases, we map each to its most frequent matching synset (according to WordNet lexeme counts). This resulted in an overall mapping accuracy of 88% and a recall of 98.5% (Fig. 5). The most common synsets extracted from region descriptions, QAs, and objects are shown in Fig. 33.

**Attributes** We canonicalize attributes from the crowd-extracted attributes present in our scene graphs. The “attribute” designation encompasses a wide range of grammatical parts of speech. Because part-of-speech taggers rely on high-level syntax information and thus fail on the disjoint elements of our scene graphs, we normalize each attribute based on morphology alone (so-called “stemming” (Bird 2006)). Then, as with objects, we map each attribute phrase to the most frequent matching WordNet synset. We include 15 hand-mapped rules to address common failure cases in which WordNet’s frequency counts prefer abstract senses of words over the spatial senses present in visual data, e.g. *short.a.01: limited in duration* over *short.a.02: lacking in length*. For verification, we randomly sample 200 attributes, produce ground-truth mappings by hand, and compare them to the results of our algorithm. This resulted in a recall of 95.9% and a mapping accuracy of 85.7%. The most common attribute synsets are shown in Fig. 34a.

**Relationships** As with attributes, we canonicalize the relationships isolated in our scene graphs. We exclude prepositions, which are not recognized in WordNet, leaving a set primarily composed of verb relationships. Since the meanings of verbs are highly dependent upon their morphology and syntactic placement (e.g. passive cases, prepositional phrases), we map the structure of each relationship to the appropriate WordNet sentence frame and only

consider those WordNet synsets with matching sentence frames. For each verb-synset pair, we then consider the root hypernym of that synset to reduce potential noise from WordNet’s fine-grained sense distinctions. We also include 20 hand-mapped rules, again to correct for WordNet’s lower representation of concrete or spatial senses; for example, the concrete *hold.v.02: have or hold in one’s hand or grip* is less frequent in WordNet than the abstract *hold.v.01: cause to continue in a certain state*. For verification, we again randomly sample 200 relationships and compare the results of our canonicalization against ground-truth mappings. This resulted in a recall of 88.5% and a mapping accuracy of 92.9%. While several datasets, such as VerbNet (Schuler 2005) and FrameNet (Baker et al. 1998), include semantic restrictions or frames to improve classification, there is no comprehensive method of mapping to those restrictions or frames. The most common relationship synsets are shown in Fig. 34b.

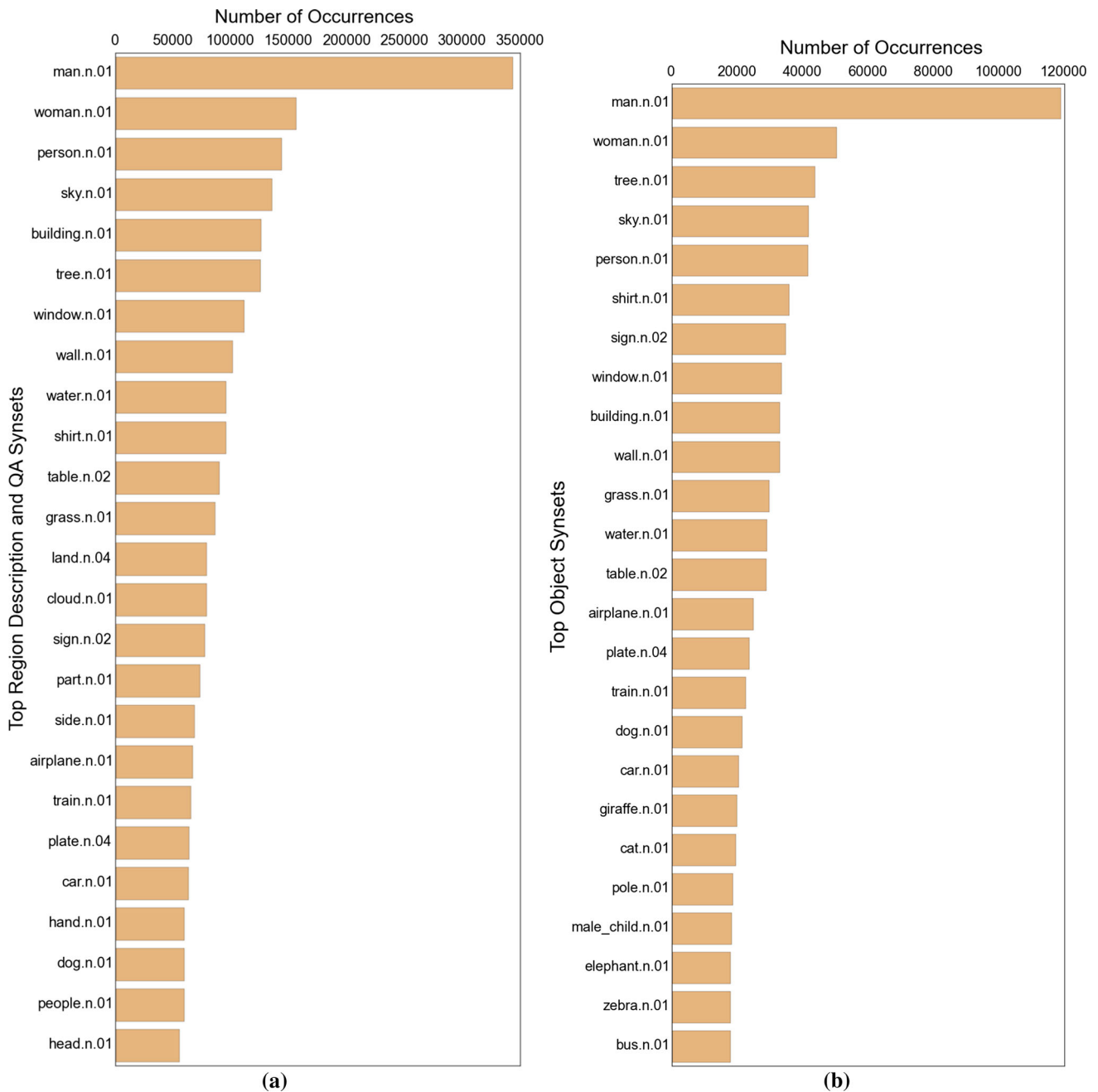
## 6 Experiments

Thus far, we have presented the Visual Genome dataset and analyzed its individual components. With such rich information provided, numerous perceptual and cognitive tasks can be tackled. In this section, we aim to provide baseline experimental results using components of Visual Genome that have not been extensively studied.

Object detection is already a well-studied problem (Everingham et al. 2010; Girshick et al. 2014; Sermanet et al. 2013; Girshick 2015; Ren et al. 2015b). Similarly, region graphs and scene graphs have been shown to improve semantic image retrieval (Johnson et al. 2015; Schuster et al. 2015). We therefore focus on the remaining components, i.e. *attributes*, *relationships*, *region descriptions*, and *question answer pairs*.

In Sect. 6.1, we present results for two experiments on attribute prediction. In the first, we treat attributes independently from objects and train a classifier for each attribute, i.e. a classifier for *red* or a classifier for *old*, as in Malisiewicz et al. (2008), Varma and Zisserman (2005), Ferrari and Zisserman (2007), Farhadi et al. (2009) and Johnson et al. (2015). In the second experiment, we learn object and attribute classifiers *jointly* and predict object-attribute pairs (e.g. predicting that an apple is red), as in Sadeghi and Farhadi (2011).

In Sect. 6.2, we present two experiments on relationship prediction. In the first, we aim to predict the predicate between two objects, e.g. predicting the predicate *kicking* or *wearing* between two objects. This experiment is synonymous with existing work in action recognition (Gupta et al. 2009; Ramanathan et al. 2015). In another experiment,

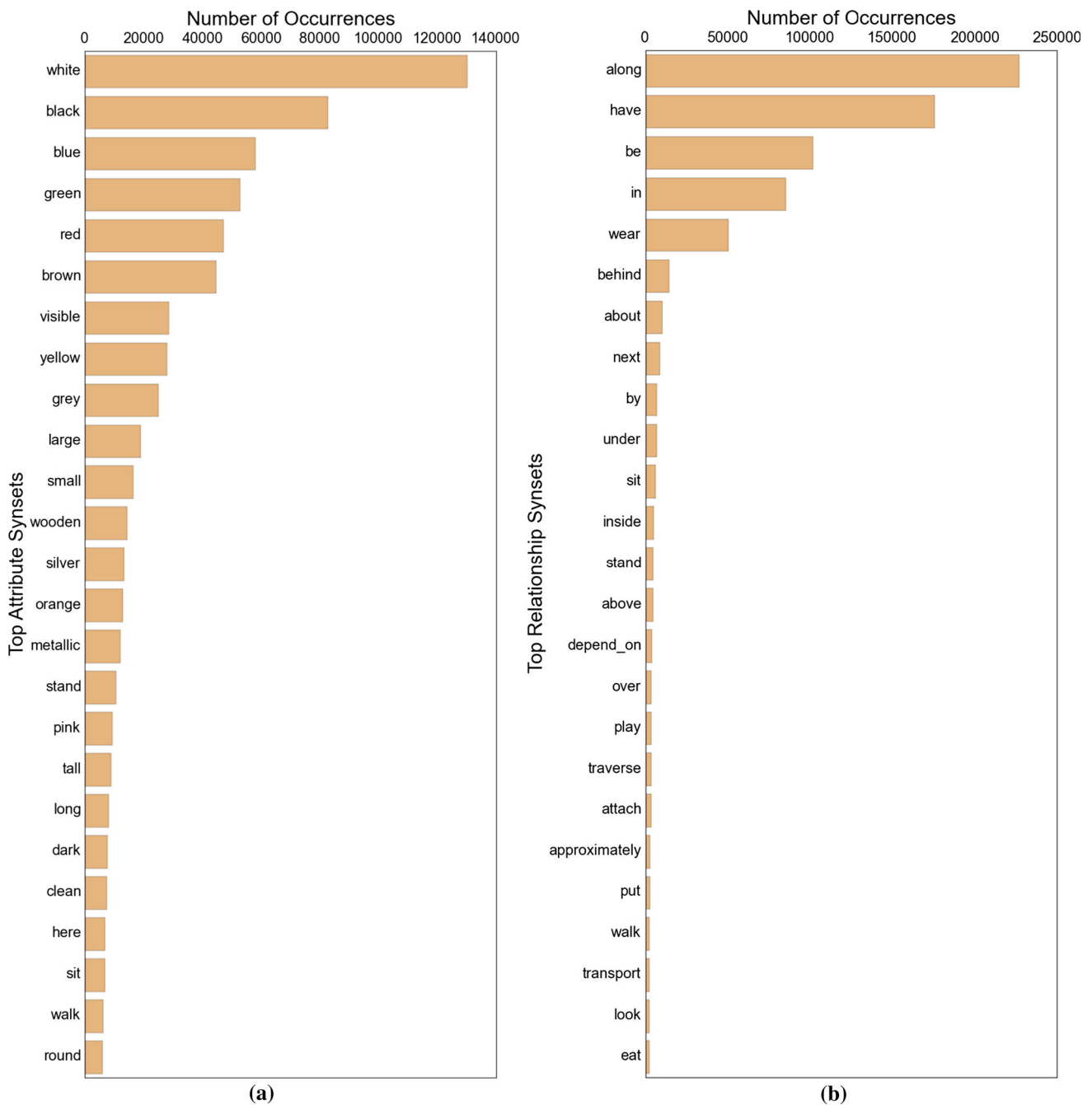


**Fig. 33** Distribution of the 25 most common synsets mapped from the words and phrases extracted from region descriptions which represent objects in **a** region descriptions and question answers and **b** objects

we study relationships by classifying jointly the objects and the predicate (e.g. predicting *kicking(man, ball)*); we show that this is a very difficult task due to the high variability in the appearance of a relationship (e.g. the *ball* might be on the ground or in mid-air above the *man*). These experiments are generalizations of tasks that study spatial relationships between objects and ones that jointly reason about the interaction of humans with objects (Yao and Fei-Fei 2010; Prest et al. 2012).

In Sect. 6.3 we present results for region captioning. This task is closely related to image captioning (Chen et al. 2015); however, results from the two are not directly comparable, as region descriptions are short, incomplete sentences. We train one of the top 16 state-of-the-art image caption generators (Karpathy and Fei-Fei 2015) on (1) our dataset to generate region descriptions and on (2) Flickr30K (Young et al. 2014) to generate sentence descriptions. To compare results between the two train-





**Fig. 34** Distribution of the 25 most common synsets mapped from **a** attributes and **b** relationships

ing approaches, we use simple templates to convert region descriptions into complete sentences. For a more robust evaluation, we validate the descriptions we generate using human judgment.

Finally, in Sect. 6.4, we experiment on visual question answering, i.e. given an image and a question, we attempt to provide an answer for the question. We report results on the retrieval of the correct answer from a list of existing answers.

## 6.1 Attribute Prediction

Attributes are becoming increasingly important in the field of computer vision, as they offer higher-level semantic cues for various problems and lead to a deeper understanding of images. We can express a wide variety of properties through attributes, such as form (*sliced*), function (*decorative*), sentiment (*angry*), and even intention (*helping*). Distinguishing between similar objects (*Isola*

et al. 2015) leads to finer-grained classification, while describing a previously unseen class through attributes shared with known classes can enable “zero-shot” learning (Farhadi et al. 2009; Lampert et al. 2009). Visual Genome is the largest dataset of attributes, with 26 attributes per image for more than 2.8 million attributes.

**Setup** For both experiments, we focus on the 100 most common attributes in our dataset. We only use objects that occur at least 100 times and are associated with one of the 100 attributes in at least one image. For both experiments, we follow a similar data pre-processing pipeline. First, we lowercase, lemmatize (Bird 2006), and strip excess whitespace from all attributes. Since the number of examples per attribute class varies, we randomly sample 500 attributes from each category (if fewer than 500 are in the class, we take all of them).

We end up with around 50,000 attribute instances and 43,000 object-attribute pair instances in total. We use 80% of the images for training and 10% each for validation and testing. Because each image has about the same number of examples, this results in an approximately 80–10–10% split over the attributes themselves. The input data for this experiment is the cropped bounding box of the object associated with each attribute.

We train an attribute predictor by using features learned from a convolutional neural network. Specifically, we use a 16-layer VGG network (Simonyan and Zisserman 2014) pre-trained on ImageNet and fine-tune it for both of these experiments using the 50,000 attribute and 43,000 object-attribute pair instances respectively. We modify the network so that the learning rate of the final fully-connected layer is 10 times that of the other layers, as this improves convergence time. Convergence is measured as the performance on the validation set. We use a base learning rate of 0.001, which we scale by 0.1 every 200 iterations, and momentum and weight decays of 0.9 and 0.0005 respectively. We use the fine-tuned features from the network and train 100 individual SVMs (Hearst et al. 1998) to predict each attribute. We output multiple attributes for each bounding box input. For the second experiment, we also output the object class.

**Results** Table 6 shows results for both experiments. For the first experiment on attribute prediction, we converge after around 700 iterations with 18.97% top-one accuracy and 43.11% top-five accuracy. Thus, attributes (like objects) are visually distinguishable from each other. For the second experiment where we also predict the object class, we converge after around 400 iterations with 43.17% top-one accuracy and 71.97% top-five accuracy. Predicting objects jointly with attributes increases the top-one accuracy from 18.97% to 43.17%. This implies that some attributes occur exclusively with a small number of objects. Additionally, by

**Table 6** (First row) Results for the attribute prediction task where we only predict attributes for a given image crop. (Second row) Attribute-object prediction experiment where we predict both the attributes as well as the object from a given crop of the image

	Top-1 accuracy (%)	Top-5 accuracy (%)
Attribute	18.97	43.11
Object-attribute	43.17	71.97

jointly learning attributes with objects, we increase the inter-class variance, making the classification process an easier task.

Figure 35a shows example predictions for the first attribute prediction experiment. In general, the model is good at associating objects with their most salient attributes, for example, animal with stuffed and elephant with grazing. However, the crowdsourced ground truth answers sometimes do not contain all valid attributes, so the model is incorrectly penalized for some accurate/true predictions. For example, the white stuffed animal is correct but evaluated as incorrect.

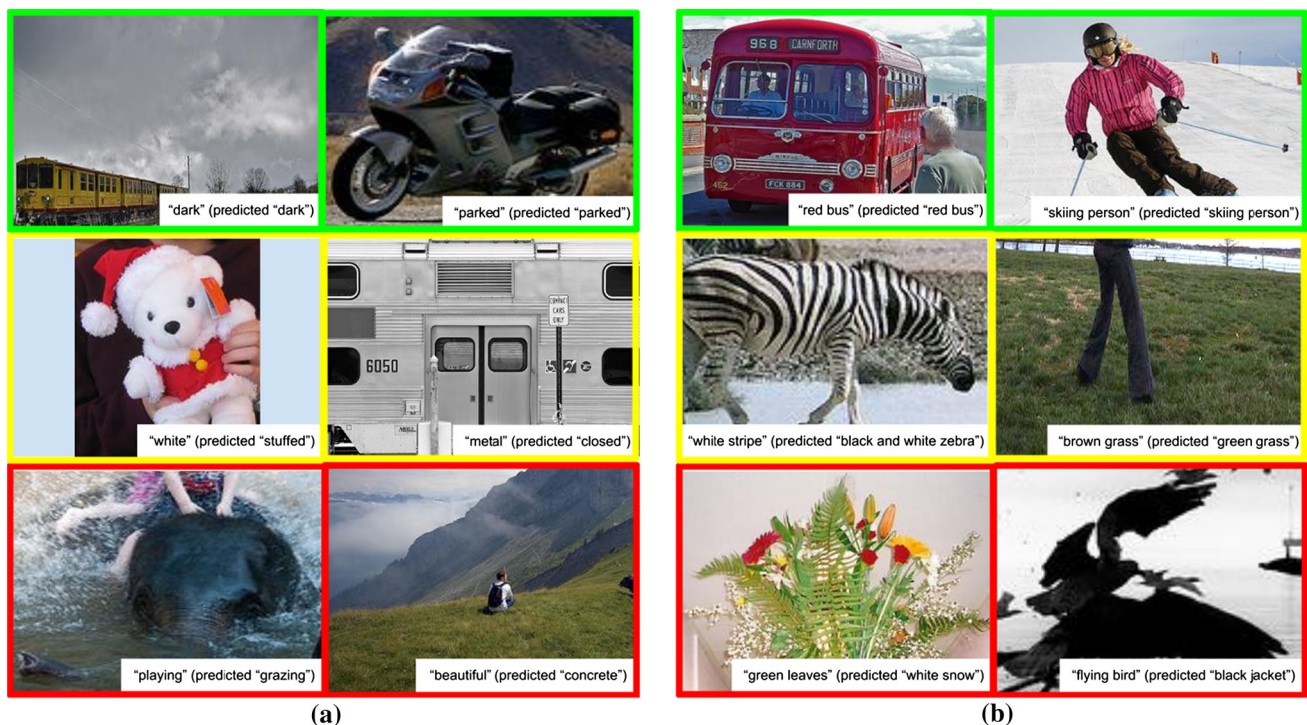
Figure 35b shows example predictions for the second experiment in which we also predict the object. While the results in the second row might be considered correct, to keep a consistent evaluation, we mark them as incorrect. For example, the predicted “green grass” might be considered subjectively correct even though it is annotated as “brown grass”. For cases where the objects are not clearly visible but are abstract outlines, our model is unable to predict attributes or objects accurately. For example, it thinks that the “flying bird” is actually a “black jacket”.

The attribute clique graphs in Sect. 5.4 clearly show that learning attributes can help us identify types of objects. This experiment strengthens that insight. We learn that studying attributes together with objects can improve attribute prediction.

## 6.2 Relationship Prediction

While objects are the core building blocks of an image, relationships put them in context. These relationships help distinguish between images that contain the same objects but have different holistic interpretations. For example, an image of “a man riding a bike” and “a man falling off a bike” both contain man and bike, but the relationship (riding vs. falling\_off) changes how we perceive both situations. Visual Genome is the largest known dataset of relationships, with more than 2.3 million relationships and an average of 21 relationships per image.

**Setup** The setups of both experiments are similar to those of the experiments we performed on attributes. We again focus



**Fig. 35** **a** Example predictions from the attribute prediction experiment. Attributes in the first row are predicted correctly, those in the second row differ from the ground truth but still correctly classify an attribute in the image, and those in the third row are classified

incorrectly. The model tends to associate objects with attributes (e.g. elephant with grazing). **b** Example predictions from the joint object-attribute prediction experiment

on the top 100 most frequent relationships. We lowercase, lemmatize (Bird 2006), and strip excess whitespace from all relationships. We end up with around 34,000 unique relationship types and 27,000 unique subject-relationship-object triples for training, validation, and testing. The input data to the experiment is the image region containing the union of the bounding boxes of the subject and object (essentially, the bounding box containing the two object boxes). We fine-tune a 16-layer VGG network (Simonyan and Zisserman 2014) with the same learning rates mentioned in Sect. 6.1.

**Results** Overall, we find that relationships are only slightly visually distinct enough for our discriminative model to learn effectively. Table 7 shows results for both experiments. For relationship classification, we converge after around 800 iterations with 8.74% top-one accuracy and 29.69% top-five accuracy. Unlike attribute prediction, the accuracy results for relationships are much lower because of the high intra-class variability of most relationships. For the second experiment jointly predicting the relationship and its two object classes, we converge after around 450 iterations with 25.83% top-one accuracy and 65.57% top-five accuracy. We notice that object classification aids relationship prediction. Some relationships occur with some objects and never others; for example, the relationship *drive* only occurs with the object *person* and never with any other objects (dog, chair, etc.).

**Table 7** Results for relationship classification (*first row*) and joint classification (*second row*) experiments

	Top-1 accuracy (%)	Top-5 accuracy (%)
Relationship	8.74	26.69
Sub./Rel./Obj.	25.83	65.57

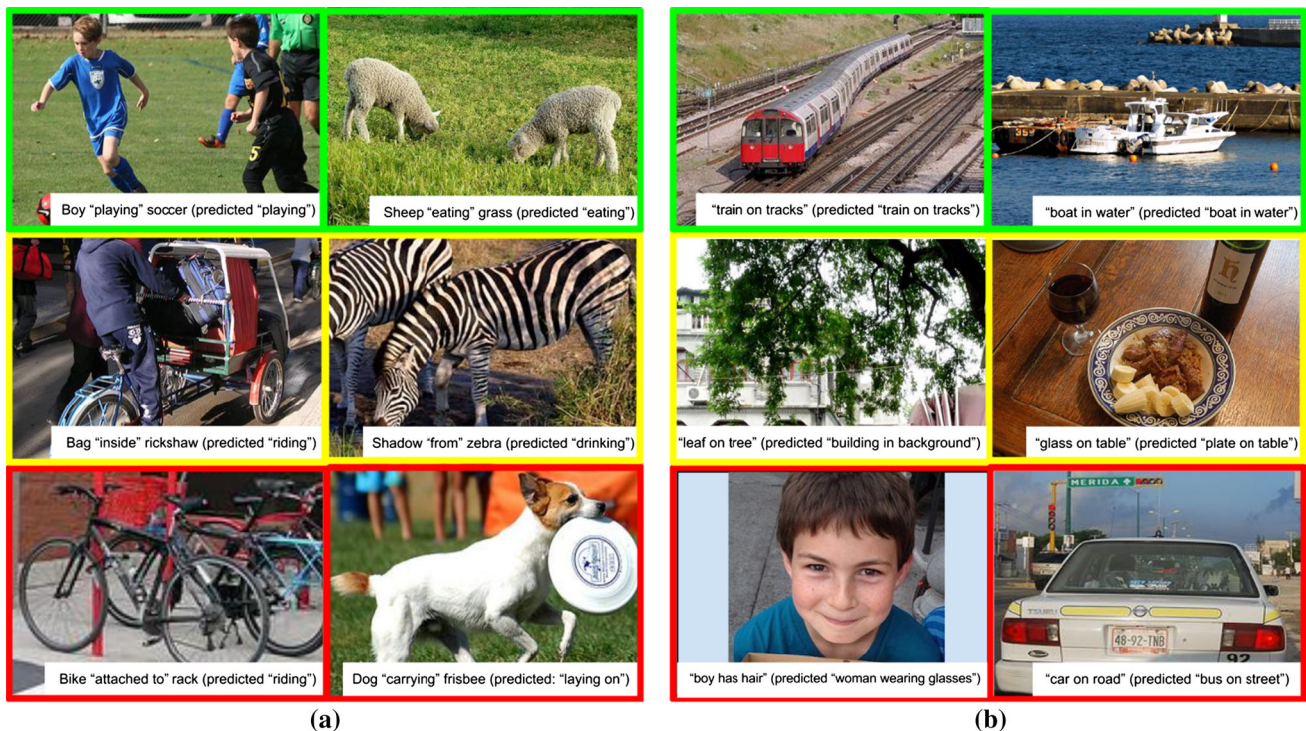
Figure 36a shows example predictions for the relationship classification experiment. In general, the model associates object categories with certain relationships (e.g. animals with eating or drinking, bikes with riding, and kids with playing).

Figure 36b, structured as in Fig. 36a, shows example predictions for the joint prediction of relationships with its objects. The model is able to predict the salient features of the image (e.g. “boat in water”) but fails to distinguish between different objects (e.g. boy vs. woman and car vs. bus in the *bottom row*).

### 6.3 Generating Region Descriptions

Generating sentence descriptions of images has gained popularity as a task in computer vision (Kiros et al. 2014; Mao et al. 2014; Karpathy and Fei-Fei 2015; Vinyals et al. 2015); however, current state-of-the-art models fail to describe all





**Fig. 36** **a** Example predictions from the relationship prediction experiment. Relationships in the first row are predicted correctly, those in the second row differ from the ground truth but still correctly classify a relationship in the image, and those in the third row are classified incorrectly. The model learns to associate animals leaning towards the ground

as eating or drinking and bikes with riding. **b** Example predictions from the relationship-objects prediction experiment. The figure is organized in the same way as **a**. The model is able to predict the salient features of the image but fails to distinguish between different objects (e.g. boy and woman and car and bus in the bottom row)

the different events captured in an image and instead provide only a high-level summary of the image. In this section, we test how well state-of-the-art models can caption the details of images. For both experiments, we use the NeuralTalk model (Karpathy and Fei-Fei 2015), since it not only provides state-of-the-art results but also is shown to be robust enough for predicting short descriptions. We train NeuralTalk on the Visual Genome dataset for region descriptions and on Flickr30K (Young et al. 2014) for full sentence descriptions. As a model trained on other datasets would generate complete sentences and would not be comparable (Chen et al. 2015) to our region descriptions, we convert all region descriptions generated by our model into complete sentences using predefined templates (Hou et al. 2002).

**Setup** For training, we begin by preprocessing region descriptions; we remove all non-alphanumeric characters and lowercase and strip excess whitespace from them. We have 5,406,939 region descriptions in total. We end up with 3,784,857 region descriptions for training – 811,040 each for validation and testing. Note that we ensure descriptions of regions from the same image are exclusively in the training, validation, or testing set. We feed the bounding boxes of the regions through the pretrained VGG 16-layer

network (Simonyan and Zisserman 2014) to get the 4096-dimensional feature vectors of each region. We then use the NeuralTalk (Karpathy and Fei-Fei 2015) model to train a long short-term memory (LSTM) network (Hochreiter and Schmidhuber 1997) to generate descriptions of regions. We use a learning rate of 0.001 trained with rmsprop (Dauphin et al. 2015). The model converges after four days.

For testing, we crop the ground-truth region bounding boxes of images and extract their 4096-dimensional 16-layer VGG network (Simonyan and Zisserman 2014) features. We then feed these vectors through the pretrained NeuralTalk model to get predictions for region descriptions.

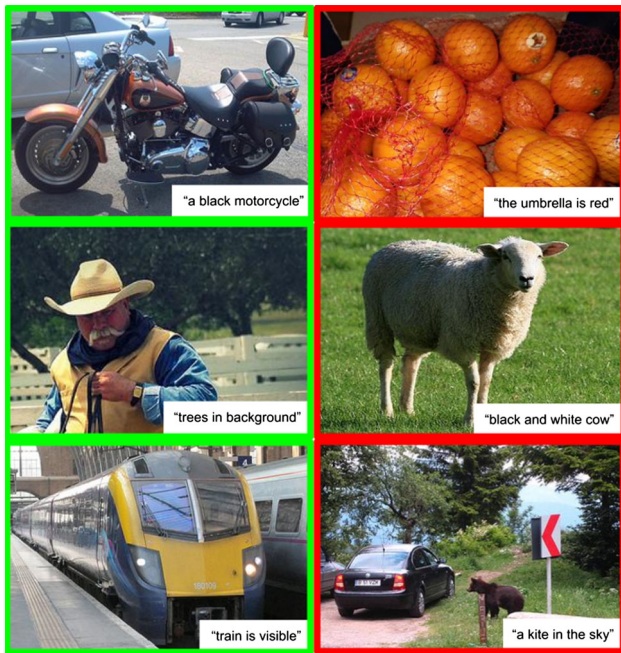
**Results** Table 8 shows the results for the experiment. We calculate BLEU (Papineni et al. 2002), CIDEr (Vedantam et al. 2015a), and METEOR (Denkowski and Lavie 2014) scores (Chen et al. 2015) between the generated descriptions and their ground-truth descriptions. In all cases, the model trained on VisualGenome performs better. Moreover, we asked crowd workers to evaluate whether a generated description was correct—we got 1.6 and 43.03% for models trained on Flickr30K and on Visual Genome, respectively. The large increase in accuracy when the model trained on our data is due to the specificity of our dataset. Our region



**Table 8** Results for the region description generation experiment

	BLEU-1	BLEU-2	BLEU-3	BLEU-4	CIDEr	METEOR	Human
Flickr8K	0.09	0.01	0.002	0.0004	0.05	0.04	1.6%
VG	0.17	0.05	0.02	0.01	0.30	0.09	43.03%

Scores in the first row are for the region descriptions generated from the NeuralTalk model trained on *Flickr8K*, and those in the second row are for those generated by the model trained on Visual Genome data. *BLEU*, *CIDEr*, and *METEOR* scores all compare the predicted description to a ground truth in different ways



**Fig. 37** Example predictions from the region description generation experiment by a model trained on Visual Genome region descriptions. Regions in the first column (*left*) accurately describe the region, and those in the second column (*right*) are incorrect and unrelated to the corresponding region

descriptions are shorter and cover a smaller image area. In comparison, the Flickr30K data are generic descriptions of entire images with multiple events happening in different regions of the image. The model trained on our data is able to make predictions that are more likely to concentrate on the specific part of the image it is looking at, instead of generating a summary description. The objectively low accuracy in both cases illustrates that current models are unable to reason about complex images.

Figure 37 shows examples of regions and their predicted descriptions. Since many examples have short descriptions, the predicted descriptions are also short as expected; however, this causes the model to fail to produce more descriptive phrases for regions with multiple objects or with distinctive objects (i.e. objects with many attributes). While we use templates to convert region descriptions into sentences, future work can explore smarter approaches to combine region descriptions and generate a paragraph connecting all the regions into one coherent description.

## 6.4 Question Answering

Visual Genome is currently the largest dataset of visual question answers with more than 1.7 million question and answer pairs. Each of our 108,077 images contains an average of 17 question answer pairs. Answering questions requires a deeper understanding of an image than generic image captioning. Question answering can involve fine-grained recognition (e.g. “What is the breed of the dog?”), object detection (e.g. “Where is the kite in the image?”), activity recognition (e.g. “What is this man doing?”), knowledge base reasoning (e.g. “Is this glass full?”), and common-sense reasoning (e.g. “What street will we be on if we turn right?”).

By leveraging the detailed annotations in the scene graphs in Visual Genome, we envision building smart models that can answer a myriad of visual questions. While we encourage the construction of smart models, in this paper, we provide some baseline results to help others compare their models.

**Setup** We split the QA pairs into a training set (60%) and a test set (40%). We ensure that all images are exclusive to either the training set or the test set. We implement a simple baseline model that relies on answer frequency. The model counts the top  $k$  most frequent answers [similar to the ImageNet challenge (Russakovsky et al. 2015)] in the training set as the predictions for all the test questions, where  $k = 100, 500$ , and  $1000$ . We let a model make  $k$  different predictions. We say the model is correct on a QA if one of the  $k$  predictions matches exactly with the ground-truth answer. We report the accuracy over all test questions. This evaluation method works well when the answers are short, especially for single-word answers. However, it causes problems when the answers are long phrases and sentences. We also report humans performance [similar to previous work (Antol et al. 2015; Yu et al. 2015)] on these questions by presenting them with the image and the question along with 10 multiple choice answers out of which one of them was the ground truth and the other 9 were randomly chosen from the dataset. Other evaluation methods require word ontologies (Malinowski and Fritz 2014).

**Results** Table 9 shows the performance of the open-ended visual question answering task. These baseline results imply the long-tail distribution of the answers. Long-tail distribu-

**Table 9** Baseline QA performances in the 6 different question types

	Top-100	Top-500	Top-1000	Human
What	0.420	0.602	0.672	0.965
Where	0.096	0.324	0.418	0.957
When	0.714	0.809	0.834	0.944
Who	0.355	0.493	0.605	0.965
Why	0.034	0.118	0.187	0.927
How	0.780	0.827	0.846	0.942
Overall	0.411	0.573	0.641	0.966

We report human evaluation as well as a baseline method that predicts the most frequently occurring answer in the dataset

tion is common in existing QA datasets as well (Antol et al. 2015; Malinowski and Fritz 2014). The top 100, 500, and 1000 most frequent answers only cover 41.1%, 57.3%, and 64.1% of the correct answers. In comparison, the corresponding sets of frequent answers in VQA (Antol et al. 2015) cover 63%, 75%, and 80% of the test set answers. The “where” and “why” questions, which tend to involve spatial and common sense reasoning, tend to have more diverse answers and hence perform poorly, with performances of 9.6 and 3.4% top-100 respectively. The top 1000 frequent answers cover only 41.8 and 18.7% of the correct answers from these two question types respectively. In comparison, humans perform extremely well in all the questions types achieving an overall accuracy of 96.6%.

## 7 Future Applications and Directions

We have analyzed the individual components of this dataset and presented experiments with baseline results for tasks such as attribute classification, relationship classification, description generation, and question answering. There are, however, more applications and experiments for which our dataset can be used. In this section, we note a few potential applications that our dataset can enable.

**Dense Image Captioning** We have seen numerous image captioning papers (Kiros et al. 2014; Mao et al. 2014; Karpathy and Fei-Fei 2015; Vinyals et al. 2015) that attempt to describe an entire image with a single caption. However, these captions do not exhaustively describe every part of the scene. A natural extension to this application, which the Visual Genome dataset enables, is the ability to create dense captioning models that describe parts of the scene.

**Visual Question Answering** While visual question answering has been studied as a standalone task (Yu et al. 2015; Ren et al. 2015a; Antol et al. 2015; Gao et al. 2015), we introduce a dataset that combines all of our question answers

with descriptions and scene graphs. Future work can build supervised models that utilize various components of Visual Genome to tackle question answering.

**Image Understanding** While we have seen a surge of image captioning (Kiros et al. 2014) and question answering (Antol et al. 2015) models, there has been little work on creating more comprehensive evaluation metrics to measure how well these models are performing. Such models are usually evaluated using BLEU, CIDEr, or METEOR and other similar metrics that do not effectively measure how well these models understand the image (Chen et al. 2015). The Visual Genome scene graphs can be used as a measurement for image understanding. Generated descriptions and answers can be matched against the ground truth scene graph of an image to evaluate its corresponding model.

**Relationship Extraction** Relationship extraction has been extensively studied in information retrieval and natural language processing (Zhou et al. 2007; GuoDong et al. 2005; Culotta and Sorensen 2004; Socher et al. 2012). Visual Genome is the first large-scale visual relationship dataset. This dataset can be used to study the extraction of visual relationships (Sadeghi et al. 2015) from images, and its interactions between objects can also be used to study action recognition (Yao and Fei-Fei 2010; Ramanathan et al. 2015) and spatial orientation between objects (Gupta et al. 2009; Prest et al. 2012).

**Semantic Image Retrieval** Previous work has already shown that scene graphs can be used to improve semantic image search (Johnson et al. 2015; Schuster et al. 2015). Further methods can be explored using our region descriptions combined with region graphs. Attention-based search methods can also be explored where the area of interest specified by a query is also localized in the retrieved images.

**Completing the Set of Annotations** While Visual Genome is the most densely annotated visual dataset for cognitive image understanding, it is still not complete. In most images, it is not feasible to collect an exhaustive set of attributes and relationships for every object or pair of objects. This raises two new research questions. In computer vision, we need to develop new evaluation metrics that do not penalize models due to a lack of a complete set of annotations. In human computer interaction, we need to design new interfaces and workflows that incentivize humans to annotate visual common sense.

## 8 Conclusion

Visual Genome provides a multi-layered understanding of pictures. It allows for a multi-perspective study of an image,

from pixel-level information like objects, to relationships that require further inference, and to even deeper cognitive tasks like question answering. It is a comprehensive dataset for training and benchmarking the next generation of computer vision models. With Visual Genome, we expect these models to develop a broader understanding of our visual world, complementing computers' capacities to detect objects with abilities to describe those objects and explain their interactions and relationships. Visual Genome is a large *formalized knowledge representation* for visual understanding and a more *complete set of descriptions and question answers* that *grounds visual concepts to language*.

**Acknowledgements** We would like to start by thanking our sponsors: Stanford Computer Science Department, Yahoo Labs, The Brown Institute for Media Innovation, Toyota, Adobe and ONR MURI. Next, we specially thank Michael Stark, Yutian Li, Frederic Ren, Sherman Leung, Michelle Guo and Gavin Mai for their contributions. We thank Carsten Rother from the University of Dresden for facilitating Oliver Groth's involvement. We also thank all the thousands of crowd workers for their diligent contribution to Visual Genome. Finally, we thank all members of the Stanford Vision Lab and Stanford HCI Group for their useful comments and discussions.

**Open Access** This article is distributed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made.

## References

- Antol, S., Agrawal, A., Lu, J., Mitchell, M., Batra, D., Zitnick, C. L., et al. (2015). VQA: Visual question answering. In *International conference on computer vision (ICCV)*.
- Antol, S., Zitnick, C. L., & Parikh, D. (2014). Zero-shot learning via visual abstraction. In *European conference on computer vision* (pp. 401–416). Springer.
- Baker, C. F., Fillmore, C. J., & Lowe, J. B. (1998). The Berkeley framenet project. In *Proceedings of the 36th annual meeting of the association for computational linguistics and 17th international conference on computational linguistics—Volume 1, ACL'98* (pp. 86–90). Stroudsburg, PA: Association for Computational Linguistics.
- Betteridge, J., Carlson, A., Hong, S. A., Hruschka, E. R. Jr., Law, E. L., Mitchell, T. M., et al. (2009). Toward never ending language learning. In *AAAI spring symposium: Learning by reading and learning to read* (pp. 1–2).
- Bird, S. (2006). NLTK: The natural language toolkit. In *Proceedings of the COLING/ACL on interactive presentation sessions* (pp. 69–72). Association for Computational Linguistics.
- Bruner, J. (1990). Culture and human development: A new look. *Human Development*, 33(6), 344–355.
- Bunescu, R. C., & Mooney, R. J. (2005). A shortest path dependency kernel for relation extraction. In *Proceedings of the conference on human language technology and empirical methods in natural language processing* (pp. 724–731). Association for Computational Linguistics.
- Chang, A. X., Savva, M., & Manning, C. D. (2014). Semantic parsing for text to 3D scene generation. In *ACL 2014* (p. 17).
- Chen, X., Fang, H., Lin, T.-Y., Vedantam, R., Gupta, S., Dollar, P., et al. (2015). *Microsoft COCO captions: Data collection and evaluation server*. [arXiv:1504.00325](https://arxiv.org/abs/1504.00325).
- Chen, X., & Lawrence Zitnick, C. (2015). Mind's eye: A recurrent visual representation for image caption generation. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 2422–2431).
- Chen, X., Liu, Z., & Sun, M. (2014). A unified model for word sense representation and disambiguation. In *EMNLP* (pp. 1025–1035). Citeseer.
- Chen, X., Shrivastava, A., & Gupta, A. (2013). Neil: Extracting visual knowledge from web data. In *2013 IEEE international conference on computer vision (ICCV)* (pp. 1409–1416). IEEE.
- Choi, W., Chao, Y.-W., Pantofaru, C., & Savarese, S. (2013). Understanding indoor scenes using 3D geometric phrases. In *2013 IEEE conference on computer vision and pattern recognition (CVPR)* (pp. 33–40). IEEE.
- Culotta, A., & Sorensen, J. (2004). Dependency tree kernels for relation extraction. In *Proceedings of the 42nd annual meeting on association for computational linguistics* (p. 423). Association for Computational Linguistics.
- Dauphin, Y., de Vries, H., & Bengio, Y. (2015). Equilibrated adaptive learning rates for non-convex optimization. In *Advances in neural information processing systems* (pp. 1504–1512).
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., & Fei-Fei, L. (2009). Imagenet: A large-scale hierarchical image database. In *IEEE conference on computer vision and pattern recognition, 2009 (CVPR 2009)* (pp. 248–255). IEEE.
- Denkowski, M., & Lavie, A. (2014). Meteor universal: Language specific translation evaluation for any target language. In *Proceedings of the ninth workshop on statistical machine translation*. Citeseer.
- Dollar, P., Wojek, C., Schiele, B., & Perona, P. (2012). Pedestrian detection: An evaluation of the state of the art. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(4), 743–761.
- Donahue, J., Anne Hendricks, L., Guadarrama, S., Rohrbach, M., Venugopalan, S., Saenko, K., et al. (2015). Long-term recurrent convolutional networks for visual recognition and description. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 2625–2634).
- Everingham, M., Van Gool, L., Williams, C. K., Winn, J., & Zisserman, A. (2010). The pascal visual object classes (VOC) challenge. *International Journal of Computer Vision*, 88(2), 303–338.
- Fang, H., Gupta, S., Iandola, F., Srivastava, R. K., Deng, L., Dollár, P., et al. (2015). From captions to visual concepts and back. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 1473–1482).
- Farhadi, A., Endres, I., Hoiem, D., & Forsyth, D. (2009). Describing objects by their attributes. In *IEEE conference on computer vision and pattern recognition, 2009 (CVPR 2009)* (pp. 1778–1785). IEEE.
- Farhadi, A., Hejrati, M., Sadeghi, M. A., Young, P., Rashtchian, C., Hockenmaier, J., et al. (2010). Every picture tells a story: Generating sentences from images. In *Computer vision—ECCV 2010* (pp. 15–29). Springer.
- Fei-Fei, L., Fergus, R., & Perona, P. (2007). Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories. *Computer Vision and Image Understanding*, 106(1), 59–70.
- Ferrari, V., & Zisserman, A. (2007). Learning visual attributes. In *Advances in neural information processing systems* (pp. 433–440).
- Ferrucci, D., Brown, E., Chu-Carroll, J., Fan, J., Gondek, D., Kalyanpur, A. A., et al. (2010). Building watson: An overview of the deepqa project. *AI Magazine*, 31(3), 59–79.
- Firestone, C., & Scholl, B. J. (2015). Cognition does not affect perception: Evaluating the evidence for top-down effects. *Behavioral and brain sciences* (pp. 1–72).



- Forbus, K. D. (1984). Qualitative process theory. *Artificial Intelligence*, 24(1), 85–168.
- Gao, H., Mao, J., Zhou, J., Huang, Z., Wang, L., & Xu, W. (2015). Are you talking to a machine? Dataset and methods for multilingual image question. In *Advances in neural information processing systems* (pp. 2296–2304).
- Geman, D., Geman, S., Hallonquist, N., & Younes, L. (2015). Visual Turing test for computer vision systems. *Proceedings of the National Academy of Sciences*, 112(12), 3618–3623.
- Girshick, R. (2015). Fast R-CNN. In *Proceedings of the IEEE international conference on computer vision* (pp. 1440–1448).
- Girshick, R., Donahue, J., Darrell, T., & Malik, J. (2014). Rich feature hierarchies for accurate object detection and semantic segmentation. In *2014 IEEE conference on computer vision and pattern recognition (CVPR)* (pp. 580–587). IEEE.
- Goering, C., Rodner, E., Freytag, A., & Denzler, J. (2014). Nonparametric part transfer for fine-grained recognition. In *2014 IEEE conference on computer vision and pattern recognition (CVPR)* (pp. 2489–2496). IEEE.
- Griffin, G., Holub, A., & Perona, P. (2007). Caltech-256 object category dataset. Technical Report 7694.
- GuoDong, Z., Jian, S., Jie, Z., & Min, Z. (2005). Exploring various knowledge in relation extraction. In *Proceedings of the 43rd annual meeting on association for computational linguistics* (pp. 427–434). Association for Computational Linguistics.
- Gupta, A., & Davis, L. S. (2008). Beyond nouns: Exploiting prepositions and comparative adjectives for learning visual classifiers. In *Computer vision—ECCV 2008* (pp. 16–29). Springer.
- Gupta, A., Kembhavi, A., & Davis, L. S. (2009). Observing human–object interactions: Using spatial and functional compatibility for recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(10), 1775–1789.
- Hayes, P. J. (1978). *The naive physics manifesto*. Geneva: Institut pour les études sémantiques et cognitives/Université de Genève.
- Hayes, P. J. (1985). The second naive physics manifesto. *Theories of the commonsense world* (pp. 1–36).
- Hearst, M. A., Dumais, S. T., Osman, E., Platt, J., & Scholkopf, B. (1998). Support vector machines. *IEEE Intelligent Systems and their Applications*, 13(4), 18–28.
- Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8), 1735–1780.
- Hodosh, M., Young, P., & Hockenmaier, J. (2013). Framing image description as a ranking task: Data, models and evaluation metrics. *Journal of Artificial Intelligence Research*, 47(1), 853–899.
- Hou, C.-S. J., Noy, N. F., & Musen, M. A. (2002). A template-based approach toward acquisition of logical sentences. In *Intelligent information processing* (pp. 77–89). Springer.
- Huang, G. B., Mattar, M., Berg, T., & Learned-Miller, E. (2008). Labeled faces in the wild: A database for studying face recognition in unconstrained environments. In *Workshop on faces in 'real-life' images: Detection, alignment, and recognition*.
- Isola, P., Lim, J. J., & Adelson, E. H. (2015). Discovering states and transformations in image collections. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 1383–1391).
- Izadinia, H., Sadeghi, F., & Farhadi, A. (2014). Incorporating scene context and object layout into appearance modeling. In *2014 IEEE conference on computer vision and pattern recognition (CVPR)* (pp. 232–239). IEEE.
- Johnson, J., Krishna, R., Stark, M., Li, L.-J., Shamma, D. A., Bernstein, M., et al. (2015). Image retrieval using scene graphs. In *IEEE conference on computer vision and pattern recognition (CVPR)*.
- Karpathy, A., & Fei-Fei, L. (2015). Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 3128–3137).
- Kiros, R., Salakhutdinov, R., & Zemel, R. (2014). Multimodal neural language models. In *Proceedings of the 31st international conference on machine learning (ICML-14)* (pp. 595–603).
- Krishna, R., Hata, K., Chen, S., Kravitz, J., Shamma, D. A., Fei-Fei, L., et al. (2016). Embracing error to enable rapid crowdsourcing. In *CHI'16-SIGCHI conference on human factors in computing system*.
- Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems* (pp. 1097–1105).
- Lampert, C. H., Nickisch, H., & Harmeling, S. (2009). Learning to detect unseen object classes by between-class attribute transfer. In *IEEE conference on computer vision and pattern recognition, 2009 (CVPR 2009)* (pp. 951–958). IEEE.
- Leacock, C., Miller, G. A., & Chodorow, M. (1998). Using corpus statistics and wordnet relations for sense identification. *Computational Linguistics*, 24(1), 147–165.
- Lebre, R., Pinheiro, P. O., & Collobert, R. (2015). *Phrase-based image captioning*. [arXiv:1502.03671](https://arxiv.org/abs/1502.03671).
- Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., et al. (2014). Microsoft COCO: Common objects in context. In *Computer vision—ECCV 2014* (pp. 740–755). Springer.
- Lu, C., Krishna, R., Bernstein, M., & Fei-Fei, L. (2016). Visual relationship detection using language priors. In *European conference on computer vision (ECCV)*. IEEE.
- Ma, L., Lu, Z., & Li, H. (2015). *Learning to answer questions from image using convolutional neural network*. [arXiv:1506.00333](https://arxiv.org/abs/1506.00333).
- Malinowski, M., & Fritz, M. (2014). A multi-world approach to question answering about real-world scenes based on uncertain input. In *Advances in neural information processing systems* (pp. 1682–1690).
- Malinowski, M., Rohrbach, M., & Fritz, M. (2015). Ask your neurons: A neural-based approach to answering questions about images. In *Proceedings of the IEEE international conference on computer vision* (pp. 1–9).
- Malisiewicz, T., Efros, A., et al. (2008). Recognition by association via learning per-exemplar distances. In *IEEE conference on computer vision and pattern recognition, 2008 (CVPR 2008)* (pp. 1–8). IEEE.
- Manning, C. D., Surdeanu, M., Bauer, J., Finkel, J., Bethard, S. J., & McClosky, D. (2014). The Stanford CoreNLP natural language processing toolkit. In *Proceedings of 52nd annual meeting of the association for computational linguistics: system demonstrations* (pp. 55–60).
- Mao, J., Xu, W., Yang, Y., Wang, J., & Yuille, A. L. (2014). *Explain images with multimodal recurrent neural networks*. [arXiv:1410.1090](https://arxiv.org/abs/1410.1090).
- Mihalcea, R., Chklovski, T. A., & Kilgariff, A. (2004). *The senseval-3 English lexical sample task*. Association for Computational Linguistics, UNT Digital Library.
- Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). *Efficient estimation of word representations in vector space*. [arXiv:1301.3781](https://arxiv.org/abs/1301.3781).
- Miller, G. A. (1995). Wordnet: a lexical database for english. *Communications of the ACM*, 38(11), 39–41.
- Niu, F., Zhang, C., Ré, C., & Shavlik, J. (2012). Elementary: Large-scale knowledge-base construction via machine learning and statistical inference. *International Journal on Semantic Web and Information Systems (IJSWIS)*, 8(3), 42–73.
- Ordonez, V., Kulkarni, G., & Berg, T. L. (2011). Im2text: Describing images using 1 million captioned photographs. In J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, & K. Weinberger (Eds.), *Advances in neural information processing systems* (Vol. 24, pp. 1143–1151). Red Hook: Curran Associates, Inc.
- Pal, A. R., & Saha, D. (2015). *Word sense disambiguation: A survey*. [arXiv:1508.01346](https://arxiv.org/abs/1508.01346).
- Papineni, K., Roukos, S., Ward, T., & Zhu, W.-J. (2002). BLEU: A method for automatic evaluation of machine translation. In



- Proceedings of the 40th annual meeting on association for computational linguistics* (pp. 311–318). Association for Computational Linguistics.
- Patterson, G., Xu, C., Su, H., & Hays, J. (2014). The sun attribute database: Beyond categories for deeper scene understanding. *International Journal of Computer Vision*, 108(1–2), 59–81.
- Perronnin, F., Sánchez, J., & Mensink, T. (2010). Improving the fisher kernel for large-scale image classification. In *Computer vision—ECCV 2010* (pp. 143–156). Springer.
- Prest, A., Schmid, C., & Ferrari, V. (2012). Weakly supervised learning of interactions between humans and objects. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(3), 601–614.
- Ramanathan, V., Li, C., Deng, J., Han, W., Li, Z., Gu, K., et al. (2015). Learning semantic relationships for better action retrieval in images. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 1100–1109).
- Ren, M., Kiros, R., & Zemel, R. (2015a). *Image question answering: A visual semantic embedding model and a new dataset*. [arXiv:1505.02074](https://arxiv.org/abs/1505.02074).
- Ren, S., He, K., Girshick, R., & Sun, J. (2015b). Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems* (pp. 91–99).
- Ronchi, M. R., & Perona, P. (2015). Describing common human visual actions in images. In X. Xie, M. W. Jones, & G. K. L. Tam (Eds.), *Proceedings of the British machine vision conference (BMVC 2015)* (pp. 52.1–52.12). BMVA Press.
- Rothe, S., & Schütze, H. (2015). *Autoextend: Extending word embeddings to embeddings for synsets and lexemes*. [arXiv:1507.01127](https://arxiv.org/abs/1507.01127).
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., et al. (2015). ImageNet large scale visual recognition challenge. *International journal of computer vision (IJCV)* (pp. 1–42).
- Russell, B. C., Torralba, A., Murphy, K. P., & Freeman, W. T. (2008). Labelme: A database and web-based tool for image annotation. *International Journal of Computer Vision*, 77(1–3), 157–173.
- Sadeghi, F., Divvala, S. K., & Farhadi, A. (2015). Viske: Visual knowledge extraction and question answering by visual verification of relation phrases. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 1456–1464).
- Sadeghi, M. A., & Farhadi, A. (2011). Recognition using visual phrases. In *2011 IEEE conference on computer vision and pattern recognition (CVPR)* (pp. 1745–1752). IEEE.
- Salehi, N., Irani, L. C., & Bernstein, M. S. (2015). We are dynamo: Overcoming stalling and friction in collective action for crowd workers. In *Proceedings of the 33rd annual ACM conference on human factors in computing systems* (pp. 1621–1630). ACM.
- Schank, R. C., & Abelson, R. P. (2013). *Scripts, plans, goals, and understanding: An inquiry into human knowledge structures*. Hove: Psychology Press.
- Schuler, K. K. (2005). *VerbNet: A broad-coverage, comprehensive verb lexicon*. Ph.D. thesis, University of Pennsylvania, Philadelphia, PA, USA (AAI3179808).
- Schuster, S., Krishna, R., Chang, A., Fei-Fei, L., & Manning, C. D. (2015). Generating semantically precise scene graphs from textual descriptions for improved image retrieval. In *Proceedings of the fourth workshop on vision and language* (pp. 70–80). Citeseer.
- Sermanet, P., Eigen, D., Zhang, X., Mathieu, M., Fergus, R., & LeCun, Y. (2013). *Overfeat: Integrated recognition, localization and detection using convolutional networks*. [arXiv:1312.6229](https://arxiv.org/abs/1312.6229).
- Silberman, N., Hoiem, D., Kohli, P., & Fergus, R. (2012). Indoor segmentation and support inference from RGBD images. In *ECCV*.
- Simonyan, K., & Zisserman, A. (2014). *Very deep convolutional networks for large-scale image recognition*. [arXiv:1409.1556](https://arxiv.org/abs/1409.1556).
- Snow, R., O'Connor, B., Jurafsky, D., & Ng, A. Y. (2008). Cheap and fast—But is it good?: Evaluating non-expert annotations for natural language tasks. In *Proceedings of the conference on empirical methods in natural language processing* (pp. 254–263). Association for Computational Linguistics.
- Socher, R., Huval, B., Manning, C. D., & Ng, A. Y. (2012). Semantic compositionality through recursive matrix-vector spaces. In *Proceedings of the 2012 joint conference on empirical methods in natural language processing and computational natural language learning* (pp. 1201–1211). Association for Computational Linguistics.
- Steinbach, M., Karypis, G., Kumar, V., et al. (2000). A comparison of document clustering techniques. In *KDD workshop on text mining*, Boston (Vol. 400, pp. 525–526).
- Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., et al. (2015). Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 1–9).
- Thomee, B., Shamma, D. A., Friedland, G., Elizalde, B., Ni, K., Poland, D., et al. (2016). YFCC100M: The new data in multimedia research. *Communications of the ACM*, 59(2), 64–73.
- Torralba, A., Fergus, R., & Freeman, W. T. (2008). 80 million tiny images: A large data set for nonparametric object and scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 30(11), 1958–1970.
- Varma, M., & Zisserman, A. (2005). A statistical approach to texture classification from single images. *International Journal of Computer Vision*, 62(1–2), 61–81.
- Vedantam, R., Lawrence Zitnick, C., & Parikh, D. (2015a). Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 4566–4575).
- Vedantam, R., Lin, X., Batra, T., Lawrence Zitnick, C., & Parikh, D. (2015b). Learning common sense through visual abstraction. In *Proceedings of the IEEE international conference on computer vision* (pp. 2542–2550).
- Vinyals, O., Toshev, A., Bengio, S., & Erhan, D. (2015). Show and tell: A neural image caption generator. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 3156–3164).
- Wah, C., Branson, S., Welinder, P., Perona, P., & Belongie, S. (2011). The Caltech-UCSD birds-200-2011 dataset. Technical Report CNS-TR-2011-001, California Institute of Technology.
- Xiao, J., Hays, J., Ehinger, K., Oliva, A., Torralba, A., et al. (2010). Sun database: Large-scale scene recognition from abbey to zoo. In *2010 IEEE conference on computer vision and pattern recognition (CVPR)* (pp. 3485–3492). IEEE.
- Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A. C., Salakhutdinov, R., Zemel, R. S., and Bengio, Y. (2015). Show, attend and tell: Neural image caption generation with visual attention. [CoRR. arXiv:1502.03044](https://arxiv.org/abs/1502.03044).
- Yang, Y., Baker, S., Kannan, A., & Ramanan, D. (2012). Recognizing proxemics in personal photos. In *2012 IEEE conference on computer vision and pattern recognition (CVPR)* (pp. 3522–3529). IEEE.
- Yao, B., & Fei-Fei, L. (2010). Modeling mutual context of object and human pose in human–object interaction activities. In *2010 IEEE conference on computer vision and pattern recognition (CVPR)* (pp. 17–24). IEEE.
- Yao, B., Yang, X., & Zhu, S.-C. (2007). Introduction to a large-scale general purpose ground truth database: methodology, annotation tool and benchmarks. In *Energy minimization methods in computer vision and pattern recognition* (pp. 169–183). Springer.
- Young, P., Lai, A., Hodosh, M., & Hockenmaier, J. (2014). From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics*, 2, 67–78.

- Yu, L., Park, E., Berg, A. C., & Berg, T. L. (2015). *Visual madlibs: Fill in the blank image generation and question answering*. [arXiv:1506.00278](#).
- Zeng, D., Liu, K., Lai, S., Zhou, G., & Zhao, J. (2014). Relation classification via convolutional deep neural network. In *Proceedings of COLING* (pp. 2335–2344).
- Zhou, G., Zhang, M., Ji, D. H., & Zhu, Q. (2007). Tree kernel-based relation extraction with context-sensitive structured parse tree information. In *EMNLP-CoNLL 2007* (p. 728).
- Zhu, J., Nie, Z., Liu, X., Zhang, B., & Wen, J.-R. (2009). Statsnowball: A statistical approach to extracting entity relationships. In *Proceedings of the 18th international conference on world wide web* (pp. 101–110). ACM.
- Zhu, Y., Fathi, A., & Fei-Fei, L. (2014). Reasoning about object affordances in a knowledge base representation. In *European conference on computer vision*.
- Zhu, Y., Zhang, C., Ré, C., & Fei-Fei, L. (2015). *Building a large-scale multimodal knowledge base system for answering visual queries*. [arXiv:1507.05670](#).
- Zitnick, C. L., & Parikh, D. (2013). Bringing semantics into focus using visual abstraction. In *2013 IEEE conference on computer vision and pattern recognition (CVPR)* (pp. 3009–3016). IEEE.