# The Open Images Dataset V4

**Unified image classification, object detection, and visual relationship detection at scale**

**Alina Kuznetsova** • **Hassan Rom** • **Neil Alldrin** • **Jasper Uijlings** • **Ivan Krasin** • **Jordi Pont-Tuset** • **Shahab Kamali** • **Stefan Popov** • **Matteo Malloci** • **Alexander Kolesnikov** • **Tom Duerig** • **Vittorio Ferrari**

Google AI

**Abstract** We present Open Images V4, a dataset of 9.2M images with unified annotations for image classification, object detection and visual relationship detection. The images have a Creative Commons Attribution license that allows to share and adapt the material, and they have been collected from Flickr without a predefined list of class names or tags, leading to natural class statistics and avoiding an initial design bias. Open Images V4 offers large scale across several dimensions: 30.1M image-level labels for 19.8k concepts, 15.4M bounding boxes for 600 object classes, and 375k visual relationship annotations involving 57 classes. For object detection in particular, we provide $15\times$ more bounding boxes than the next largest datasets (15.4M boxes on 1.9M images). The images often show complex scenes with several objects (8 annotated objects per image on average). We annotated visual relationships between them, which support visual relationship detection, an emerging task that requires structured reasoning. We provide in-depth comprehensive statistics about the dataset, we validate the quality of the annotations, we study how the performance of several modern models evolves with increasing amounts of training data, and we demonstrate two applications made possible by having unified annotations of multiple types coexisting in the same images. We hope that the scale, quality, and variety of Open Images V4 will foster further research and innovation even beyond the areas of image classification, object detection, and visual relationship detection.

**Keywords** Ground-truth dataset · Image classification · Object detection · Visual relationship detection

## 1 Introduction

Deep learning is revolutionizing many areas of computer vision. Since its explosive irruption in the ImageNet chal-

lenge (Russakovsky et al., 2015) in 2012, performance of models has been improving at an unparalleled speed. At the core of their success, however, lies the need of gargantuan amounts of annotated data to learn from. Larger and richer annotated datasets are a boon for leading-edge research in computer vision to enable the next generation of state-of-the-art algorithms.

Data is playing an especially critical role in enabling computers to interpret images as compositions of objects, an achievement that humans can do effortlessly while it has been elusive for machines so far. In particular, one would like machines to automatically identify what objects are present in the image (*image classification*), where are they precisely located (*object detection*), and which of them are interacting and how (*visual relationship detection*).

This paper presents the Open Images Dataset V4, which contains images and ground-truth annotations for the three tasks above (Figure 1). Open Images V4 has several attractive characteristics, compared to previously available datasets in these areas (Krizhevsky, 2009; Fei-Fei et al., 2006; Griffin et al., 2007; Deng et al., 2009; Russakovsky et al., 2015; Everingham et al., 2012; Gupta and Malik, 2015; Krishna et al., 2017). The images were collected from Flickr[1] without a predefined list of class names or tags, leading to natural class statistics and avoiding the initial design bias on what should be in the dataset. They were released by the authors under a *Creative Commons Attribution* (CC-BY) license that allows to share and adapt the material, even commercially; particularly so for models trained on these data, since it makes them more easily usable in any context. Also, we removed those images that appear elsewhere in the internet to reduce bias towards web image search engines, favoring complex images containing several objects. Complex images open the door to visual relationship detection, an emerging topic at the frontier of computer vision that requires structured rea-

---

[1] Image hosting service (`flickr.com`)

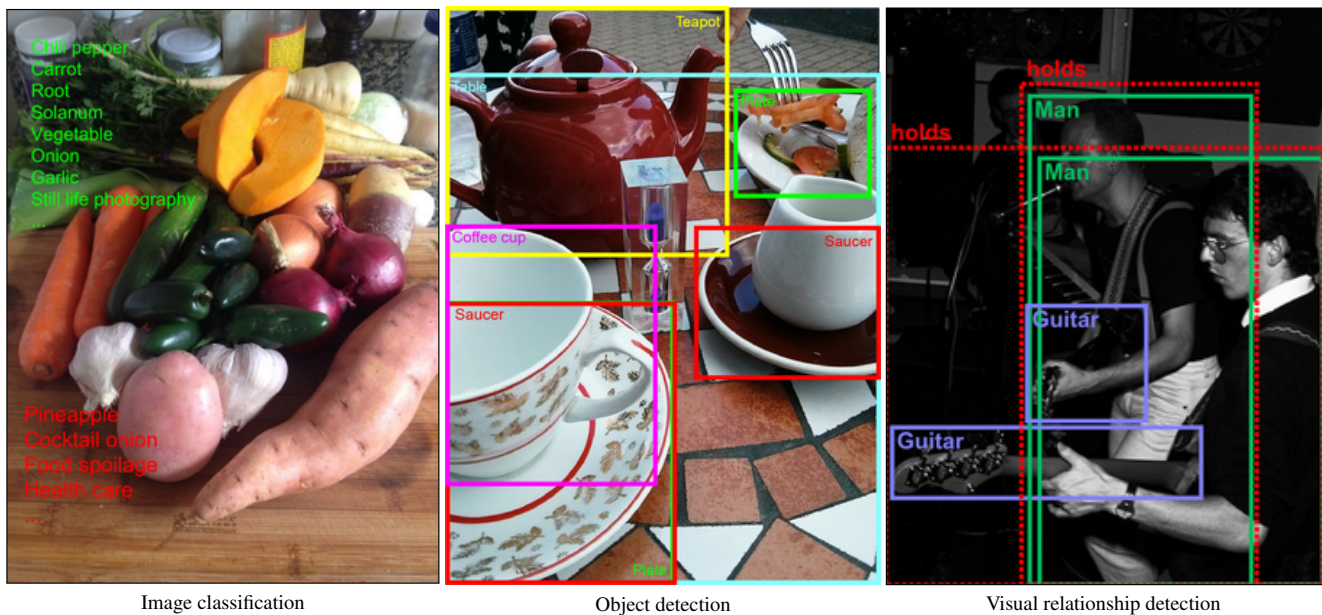| Image classification | Object detection | Visual relationship detection |

**Fig. 1 Example annotations in Open Images** for image classification, object detection, and visual relationship detection. For image classification, positive labels (present in the image) are in green while negative labels (not present in the image) are in red. For visual relationship detection, the box with a dashed contour groups the two objects that hold a certain visual relationship.

soning of the contents of an image. Section 2 further explains all the specifics about how we collected and annotated Open Images.

Open Images V4 is large scale in terms of images (9,178,275), annotations (30,113,078 image-level labels, 15,440,132 bounding boxes, 374,768 visual relationship triplets) and the number of visual concepts (*classes*) (19,794 for image-level labels and 600 for bounding boxes). This makes it ideal for pushing the limits of the data-hungry methods that dominate the state of the art. For object detection in particular, the scale of the annotations is unprecedented (15.4 million bounding boxes for 600 categories on 1.9 million images). The number of bounding boxes we provide is more than $15\times$ greater than the next largest datasets (COCO and ImageNet). Also, there are 8 annotated bounding boxes per image on average, demonstrating the complexity of the images and the richness of our annotations. We hope this will stimulate research into more sophisticated detection models that will exceed current state-of-the-art performance and will enable assessing more precisely in which situations different detectors work best. Section 3 provides an in-depth comprehensive set of statistics about Open Images V4 and compare them to previous datasets. Finally, Open Images V4 goes beyond previous datasets also in that it is *unified*: the annotations for image classification, object detection, and visual relationship detection all coexist in the same set of images. This allows for cross-task training and analysis, potentially supporting deeper insights about each of the three tasks, enabling tasks that require multiple annotation types, and stimulating progress towards genuine scene understanding.

To validate the quality of the annotations, in Section 4 we study the geometric accuracy of the bounding boxes and the recall of the image-level annotations by comparing them to annotations done by two experts and by comparing the annotators' consistency. In Section 5 we analyze the performance of several modern models for image classification and object detection, studying how their performance evolves with increasing amounts of training data, and we also report several baselines for visual relationship detection. Finally, to demonstrate the value of having unified annotations, we report in Section 6 two experiments that are made possible by them (fine-grained object detection without fine-grained box labels, and zero-shot visual relationship detection).

All the annotations, up-to-date news, box visualization tools, etc. are available on the Open Images website: `https://g.co/dataset/openimages/`. This is the first paper about Open Images, there is no previous conference or journal version.

## 2 Dataset Acquisition and Annotation

This section explains how we collected the images in the Open Images Dataset (Sec. 2.1), which classes we labeled (Sec. 2.2), and how we annotated (i) image-level labels (Sec. 2.3), (ii) bounding boxes (Sec. 2.4), and (iii) visual relationships (Sec. 2.5).

## 2.1 Image Acquisition

Images are the foundation of any good vision dataset. The Open Images Dataset differs in three key ways from most other datasets. First, all images have Creative Commons Attribution (CC-BY) license and can therefore be more easily used, with proper attribution (e.g. in commercial applications, or for crowdsourcing). Second, the images are collected starting from Flickr and then removing images that appear elsewhere on the internet. This removes simple images that appear in search engines such as Google Image Search, and therefore the dataset contains a high proportion of interesting, complex images with several objects. Third, the images are not scraped based on a predefined list of class names or tags, leading to natural class statistics and avoiding the initial design bias on what should be in the dataset.

The ~9 million images in the Open Images Dataset were collected using the following procedure:

1. Identify all Flickr images with CC-BY license. This was done in November 2015.
2. Download the original version[2] of these images and generate a copy at two resolutions:

   – *1600HQ*: Images have at most 1,600 pixels on their longest side and 1,200 pixels on their shortest. JPEG quality of 90.
   – *300K*: Images have roughly 300,000 pixels. JPEG quality of 72.

3. Extract relevant metadata of all images to give proper attribution:

   – *OriginalURL*: Flickr direct original image url.
   – *OriginalLandingURL*: Flickr image landing page.
   – *License*: Image license, a subtype of CC-BY.
   – *Author*: Flickr name of the author of the photo.
   – *Title*: Title given by the author in Flickr.
   – *AuthorProfileURL*: Link to the Flickr author profile.
   – *OriginalMD5*: MD5 hash of the original JPEG-encoded image.

4. Remove images containing inappropriate content (porn, medical, violence, memes, etc.) using the safety filters on Flickr and Google SafeSearch.
5. Remove near-duplicate images, based on low-level visual similarity.
6. Remove images that appear elsewhere on the internet. This was done for two reasons: to prevent invalid CC-BY attribution and to reduce bias towards web image search engines.



**Fig. 2 Most-frequent image-level classes**. Word size is proportional to the class counts in the train set.



**Fig. 3 Infrequent image-level classes**. Word size is inversely proportional to the class counts in the train set.

7. Recover the user-intended image orientation by comparing each original downloaded image to one of the Flickr resizes.[3]
8. Partition the images into train (9,011,219 images), validation (41,620) and test (125,436) splits (Tab. 2).

## 2.2 Classes

The set of classes included in the Open Images Dataset is derived from JFT, an internal dataset at Google with millions of images and thousands of classes (Hinton et al., 2014; Chollet, 2017; Sun et al., 2017). We selected 19,794 classes from JFT, spanning a very wide range of concepts, which serve as the image-level classes in the Open Images Dataset:

– Coarse-grained object classes (e.g. `animal`).
– Fine-grained object classes (e.g. `Pembroke welsh corgi`).
– Scene classes (e.g. `sunset` and `love`).
– Events (e.g. `birthday`).
– Materials and attributes (e.g. `leather` and `red`).

---

[2] In Flickr terms, images are served at different sizes (Thumbnail, Large, Medium, etc.). The Original size is a pristine copy of the image that was uploaded by the author.
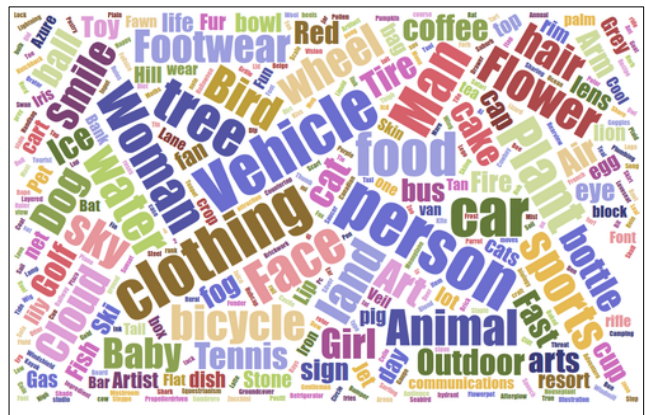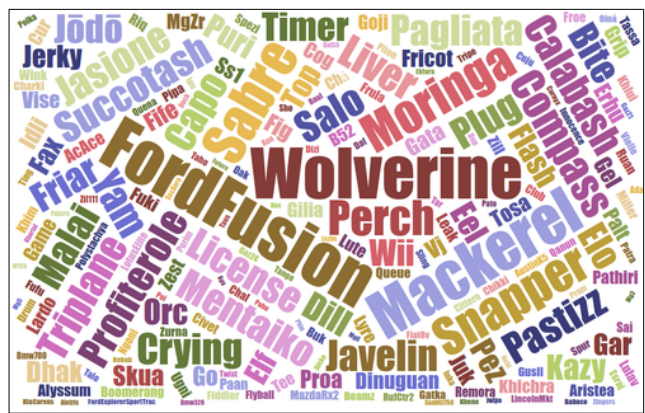
[3] More details at https://storage.googleapis.com/openimages/web/2018-05-17-rotation-information.html.

**Fig. 4 Examples of image-level labels**. Positive (green) and negative (red) image-level labels.

An overview of the most frequent and infrequent classes is shown in Figures 2 and 3.

Out of the image-level classes, we selected 600 object classes we deemed important and with a clearly defined spatial extent as *boxable*: these are classes for which we collect bounding box annotations (Sec. 2.4). A broad range of object classes are covered including animals, clothing, vehicles, food, people, buildings, sports equipment, furniture, and kitchenware. The boxable classes additionally form a hierarchy, shown in Figure 5. Figure 19 shows two example images with a wide variety of boxable classes present.

## 2.3 Image-Level Labels

Manually labeling a large number of images with the presence or absence of 19,794 different classes is not feasible not only because of the amount of time one would need, but also because of the difficulty for a human to learn and remember that many classes. In order to overcome this, we apply a computer-assisted protocol. We first apply an image classifier to generate candidate labels for all images (Sec. 2.3.1 and 2.3.2), and then ask humans to verify them (Sec. 2.3.3).

For each image, this process results in several positive (the class is present) and negative (the class is absent) labels. The presence of any other label (which has not been verified) is unknown. The negative labels are therefore valuable, as they enable to properly train discriminative classifiers even in our incomplete annotation setting. We will investigate this further in Section 5.1. Examples of positive and negative image-level labels are shown in Figure 4.

### 2.3.1 Candidate labels for test and validation

For test and validation we generate predictions for each of the 19,794 classes using a google-internal variant of the Inception V2-based image classifier (Szegedy et al., 2016), which is publicly available through the Google Cloud Vision API. This model is trained on the JFT dataset – an internal Google dataset with more than 300 million images with noisy labels (Hinton et al., 2014; Chollet, 2017). We applied this

model to the *300K* resized images in the test and validation splits. For each image, we retain all labels with a confidence score above 0.5 as candidates. We then ask humans to verify these candidate labels (Sec. 2.3.3).

### 2.3.2 Candidate labels for train

For the train split, we generate predictions by applying dozens of image classifiers. To do this, we trained various image classification models on the JFT dataset. The classification models are Google-internal and use a variety of architectures such as Inception and ResNet families. We applied all models to each of the *300K* resized images in the train split. These model predictions were used to select candidate labels to be verified by humans through stratified sampling, as explained next.

For each model we take the predictions for each image for all classes and distribute them in strata according to percentiles of their score. We then sample a certain amount of images from each class and strata to verify. The rationale behind this strategy is to have all ranges of classification scores represented in the verified sample.

Formally, for each class $c$, image $i$ is assigned to strata according to the logit scores output by the classifier $\mathbf{m}$ as:

$$\text{stratum}(i, c; \mathbf{m}) = \left\lfloor \text{logit}(i, c; \mathbf{m}) \cdot \frac{1}{w} \right\rfloor \qquad (1)$$

where $w$ is the stratum width, $\text{logit}(i, c; \mathbf{m})$ is the logit score for class $c$ from model $\mathbf{m}$ applied to image $i$, and $\lfloor \cdot \rfloor$ is the *floor* operator (i.e. rounding down to the nearest integer). Within each stratum, we sample $k$ images to be verified.

Since we perform this process for multiple classification models, the sampling of images within each stratum is not done randomly, but by selecting the $k$ images with lowest image id[4]. This way, the overall process results in far fewer than $m \cdot k$ verifications since there is high overlap of sampled image ids between models. Moreover, it encourages verifying multiple different classes on the same images: the low image ids will have high probability to be sampled for many classes, while high image ids will only be sampled for rare classes with higher-confidence model predictions.

This sampling strategy yields a good variety of examples: high confidence strata lead to a mix of easy positives and hard negatives, while low confidence strata lead to a mix of hard positives and easy negatives. We repeated this procedure for dozens of classifier models $\mathbf{m}$ using $w = 2$ and $k = 10$.[5]

Additionally, to obtain denser annotations for the 600 boxable classes, we repeated the approach in Section 2.3.2

---

[4] Image ids are generated based on hashes of the data so effectively the sampling within a stratum is pseudo-random and deterministic.

[5] Note that while in theory logit scores are unbounded, we rarely observe values outside of $[-8, 8]$ so the number of strata is bounded in practice.
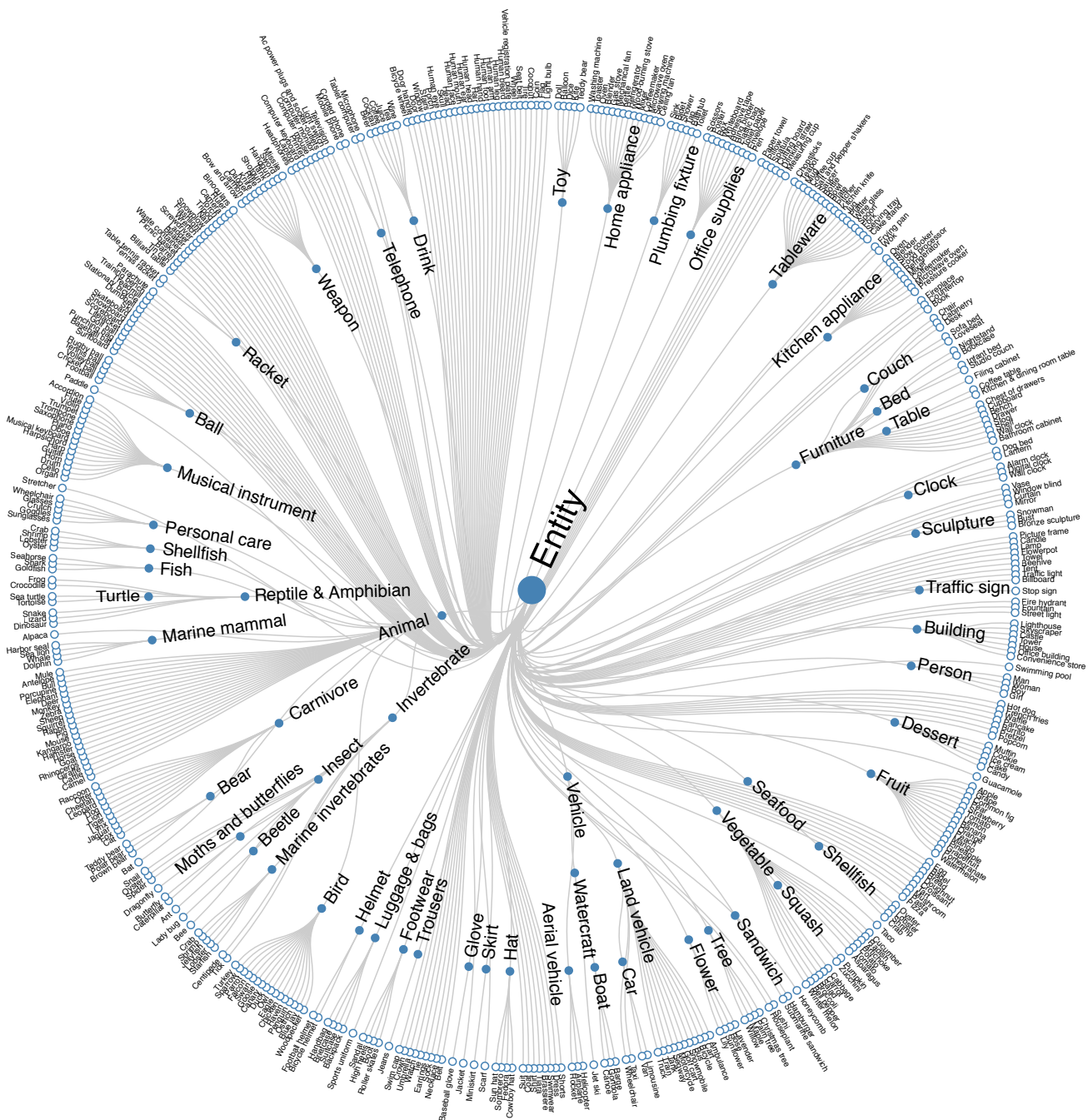
**Fig. 5 The boxable class hierarchy**. Parent nodes represent more generic concepts than their children.

on the 1.74 million training images where we annotated bounding boxes (Sec. 2.4). This generates a denser set of candidate labels for the boxable classes, to which we want to give stronger emphasis.

### 2.3.3 Human verification of candidate labels

We presented each candidate label with its corresponding image to a human annotator, who *verifies* whether the class is present in that image or not. We use two pools of annotators for such verification questions: a Google-internal pool and a crowdsource external pool. Annotators in the Google-internal pool are trained and we can provide them with extensive guidance on how to interpret and verify the presence of classes in images. The latter are Internet users that provide verifications through a crowdsourcing platform over which we cannot provide such training.

For each verification task, we use majority voting over multiple annotators. We varied the number of annotators depending on the annotator pool and the difficulty of the class

(which depends on how objective and clearly defined it is). More precisely, we used the majority of 7 annotators for crowdsource annotators. For the internal pool, we used 3 annotators for difficult classes (e.g. `ginger beer`) and 1 annotator for easy classes (including the boxable ones).

## 2.4 Bounding Boxes

We annotated bounding boxes for the 600 boxable object classes (Sec. 2.2). In this section, we first describe the guidelines which we used to define what a good bounding box is on an object (Sec. 2.4.1). Then we describe the two annotation techniques which we used (Secs. 2.4.2 and 2.4.3), followed by hierarchical de-duplication (Sec. 2.4.4), and attribute annotation (Sec. 2.4.5).

### 2.4.1 What is a perfect bounding box?

As instruction, our annotators were given the following general definition: Given a target class, a perfect box is the smallest possible box that contains all visible parts of the object (Figure 6 left). While this definition seems simple enough at first sight, there are quite a few class-dependent corner cases such as: *are straps part of a camera?*, *is water part of a fountain?*. Additionally, we found unexpected cultural differences, such as a `human hand` including the complete human arm in some parts of the world (Figure 6 right). To ensure different annotators would consistently mark the same spatial extent, we manually annotated a perfect bounding box on two examples for each of the 600 object classes. Additionally, for 20% of the classes we identified common mistakes in pilot studies. Annotators always worked on a single class at a time, and were shown the positive examples and common mistakes directly before starting each annotation session (Figure 6). This helps achieving high quality and consistency.
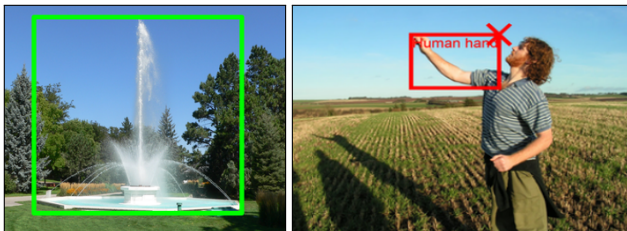


**Fig. 6 Example boxes shown to annotators**. The left shows a perfect box for `fountain`. The right shows a common mistake for `human hand`, caused by cultural differences. Only examples of the target class were shown before annotating that class.

Sometimes object instances are too close to each other to put individual boxes on them. Therefore, we also allowed annotators to draw one box around five or more heavily overlapping instances (Fig. 16 right) and mark that box with the `GroupOf` attribute (Sec. 2.4.5).

### 2.4.2 Extreme Clicking

We annotated 90% of all bounding boxes using extreme clicking, a fast box drawing technique introduced in (Papadopoulos et al., 2017). The traditional method of drawing a bounding box (Su et al., 2012), used to annotate ILSVRC (Russakovsky et al., 2015), involves clicking on imaginary corners of a tight box around the object. This is difficult as these corners are often outside the actual object and several adjustments are required to obtain a tight box. In extreme clicking, annotators are asked to click on four physical points on the object: the top, bottom, left- and right-most points. This task is more natural and these points are easy to find.

*Training annotators.* We use Google-internal annotators for drawing all boxes on the Open Images Dataset. We found it crucial to train annotators using an automated process, in the spirit of (Papadopoulos et al., 2017). Our training consists of two parts. Part one is meant to teach extreme clicking. Here annotators draw boxes on 10 objects for each of the 20 PASCAL VOC classes (Everingham et al., 2015). After each class we automatically provide feedback on which boxes were correctly or incorrectly drawn and why, by showing valid possible positions of the extreme points (Fig. 7). Part two is a qualification task in which the annotators practice both speed and accuracy. They are asked to draw 800 boxes and pass if their intersection-over-union (IoU) with the ground truth is higher than 0.84 and drawing time per box is 20 seconds or less. This sets high-quality standards, as the human expert agreement is 0.88 (Papadopoulos et al., 2017).
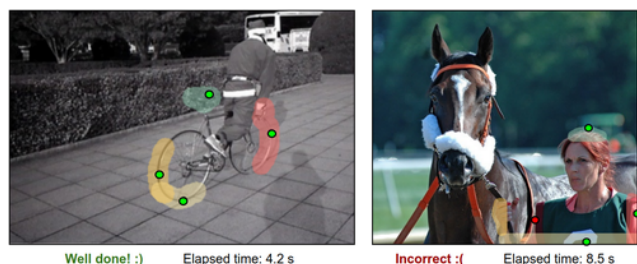


**Fig. 7 Feedback during extreme clicking training**. Left: the annotator correctly annotated a box for `bicycle`. Right: the annotator incorrectly annotated `person` (wrong point shown in red). In both cases we display the valid area for each extreme point.

*Annotation time.* On average over the complete dataset, it took 7.4 seconds to draw a single bounding box. This is much faster than the median time of 42 seconds reported for ILSVRC (Russakovsky et al., 2015; Su et al., 2012),

broken down into 25.5 seconds for drawing a box, 9.0 seconds for verifying its quality, and 7.8 seconds for checking if other instances needed to be annotated in the same image. Because of our automated training and qualification stages, we found it unnecessary to verify whether a box was drawn correctly (Sec. 4.1 for a quality analysis). Furthermore, annotators were asked to draw boxes around *all* instances of a single class in an image consecutively, removing the separate task of checking if other instances needed to be annotated. Finally, extreme clicking significantly reduced the box drawing time itself from 25.5 to 7.4 seconds.

### 2.4.3 Box Verification Series

About 10% of the bounding boxes in the training set were annotated using box verification series, in which annotators verify bounding boxes produced automatically by a learning algorithm (Papadopoulos et al., 2016). Given image-level labels, this scheme iterates between retraining the detector, re-localizing objects in the training images, and having human annotators verify bounding boxes. The verification signal is used in two ways. First, the detector is periodically retrained using all bounding boxes accepted so far, making it stronger. Second, the rejected bounding boxes are used to reduce the search space of possible object locations in subsequent iterations. Since a rejected box has an IoU $< t$ with the true bounding box, we can eliminate all candidate boxes with an IoU $\geq t$ with it ($t$ is the acceptance threshold). This is guaranteed not to remove the correct box. This strategy is effective because it eliminates those areas of the search space that matter: high scoring locations which are unlikely to contain the object.
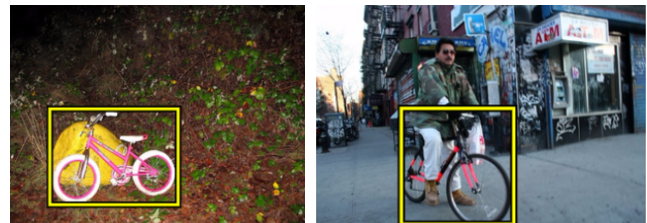
We adapted this general scheme to our operation in several ways. We make up to four attempts to obtain a bounding box for a specific class in an image. We set a higher quality criterion: we instruct annotators to accept a box if its IoU with an imaginary perfect box is greater than $t = 0.7$ (instead of $0.5$ in (Papadopoulos et al., 2016)). Additionally, to more efficiently use annotation time we did not verify boxes with a confidence score lower than $0.01$. As detector we used Faster-RCNN (Ren et al., 2015) based on Inception-ResNet (Szegedy et al., 2017) using the implementation of (Huang et al., 2017). We train our initial detector using the weakly-supervised technique with knowledge transfer described in (Uijlings et al., 2018). This uses image-level labels on the Open Images Dataset and the ILSVRC detection 2013 training set (Russakovsky et al., 2015). We retrained our detector several times during the annotation process based on all boxes accepted until that point in time. Interestingly, the final detector was truly stronger than the initial one. The annotators accepted 48% of the boxes the initial detector proposed (considering the highest-scored box for an image, if its score is $> 0.01$). This increased to 70%

for the final detector. A typical box verification series is shown in Figure 8.



**Fig. 8 Example of a box verification series for `guitar`.** The highest scored `guitar` box is shown to the annotator, who rejects it. Then the system proposes a second box, which the annotator rejects as well. Finally, the third proposed box is accepted and the process is completed.

*Training annotators.* As for extreme clicking, we found it crucial to train the annotators using an automated process. We performed several training rounds where the annotators verified 360 boxes on PASCAL VOC 2012. We automatically generated these boxes and calculated their IoU with respect to the ground truth. Ideally, the annotator should accept all boxes with IoU $> 0.7$ and reject the rest. In practice, we ignored responses on borderline boxes with IoU $\in (0.6, 0.8)$, as these are too difficult to verify. This helped relaxing the annotators, who could then focus on the important intervals of the IoU range ($[0.0, 0.6]$ and $[0.8, 1.0]$). To make training effective, after every 9 examples we provided feedback on which boxes were correctly or incorrectly verified and why (Figure 9).



This is clearly a correct box according to our definition.
You said **No**, but the correct answer is **Yes**.
The real overlap of this box is **0.96**.

Well done!
You said **Yes**, and your answer was correct.
The real overlap of this box is **0.92**.

**Fig. 9 Example feedback during the training phase of box verification series.** The target class is `bicycle`.

Figure 10 demonstrates the importance of training. It plots the acceptance rate versus the IoU for the first (———) and third (———) training rounds, and compares it to the ideal behavior (- - -). In the first round, performance was not great: 15% of boxes with almost no overlap with the object were accepted, while only 95% of boxes with very high overlap were accepted (IoU $\in [0.8, 0.9]$). Additionally, relatively poor poxes with IoU $\in [0.3, 0.6]$ were accepted 5% to 15% of the time. In contrast, after three training rounds the acceptance rate of IoU $\in [0.3, 0.6]$ was nicely below 4%, while high overlap boxes (IoU $\in [0.8, 0.9]$) were accepted 98% of

the time. 80% of the annotators were deemed qualified after three rounds of training. The other 20% needed one extra training round to reach good quality.
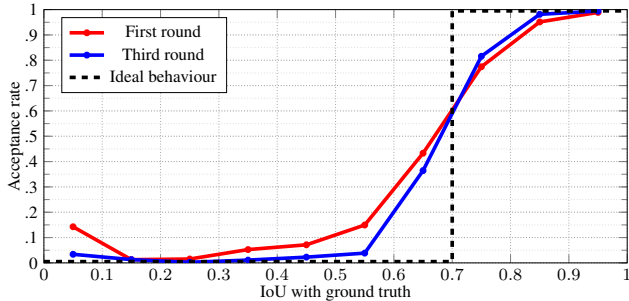


**Fig. 10 IoU versus acceptance rate for the box verification training task**. Overall, annotators do much better after three training rounds with feedback.

*Annotation time.* After a short period of time where annotators were getting used to the task, verifying a single box took 3.5 seconds on average. By dividing the total number of accepted boxes by the total time spent on verifying boxes (including cases for which box verification series failed to produce a box), we measured an average time of 8.5 seconds per box produced. This is much faster than the original annotation time for the ILSVRC boxes (25.5 seconds for manually drawing it, plus additional time to verify it, Sec. 2.4.2).

*Historical process.* We initially annotated 1.5 million boxes in the training set with box verification series. Afterwards, we co-invented extreme clicking (Papadopoulos et al., 2017). Since extreme clicking takes about the same annotation time, but it is easier to deploy and delivers more accurate boxes (Sec. 4), we used it to annotate all remaining boxes (i.e. 13.1 million in the training set, and the whole validation and test sets). Please note that in this second stage we asked annotators to draw all missing boxes for all available positive image-level labels in all images. Hence, the final dataset has a box on an object even if box verification series failed (i.e. after 4 rejected boxes).

### 2.4.4 Hierarchical de-duplication

We annotated bounding boxes for each positively verified image-level label. To prevent drawing two bounding boxes on the same object with two labels (e.g. `animal` and `zebra`), we performed hierarchical de-duplication. On the train set, before the box annotation process started, we removed all parents of another label already present in the set of image-level labels for a particular image. For example, if an image had labels `animal`, `zebra`, `car`; we annotated boxes for

`zebra` and `car`. On the validation and test splits we used a stricter, and more expensive, protocol. We first asked annotators to draw all boxes for all available labels on the image. Then we only removed a parent box (e.g. `animal`) if it overlapped with a box of a child class (e.g. `zebra`) by IoU $> 0.8$.

### 2.4.5 Attributes

We asked annotators to mark the following attributes if applicable:

`GroupOf`: the box covers more than 5 instances of the same class which heavily occlude each other.
`Partially occluded`: the object is occluded by another object in the image.
`Truncated`: the object extends outside of the image.
`Depiction`: the object is a depiction such as a cartoon, drawing, or statue.
`Inside`: the box captures the inside of an object (e.g. inside of an `aeroplane` or a `car`).

Additionally, we marked whether boxes were obtained through box verification series (Sec. 2.4.3) or through extreme clicking (Sec. 2.4.2).

`Truncated` and `Occluded` were also marked in PASCAL (Everingham et al., 2015). The purpose of `GroupOf` is similar to `crowd` in COCO (Lin et al., 2014), but its definition is different. In COCO, after having individually segmented 10-15 instances in an image, other instances in the same image were grouped together in a single, possibly disconnected, `crowd` segment.

## 2.5 Visual relationships

The Open Images Dataset is rich in terms of the of number of classes and diversity of scenes, which motivated us to annotate *visual relationships*. This will support research in the emerging topics of visual relationship detection (Lu et al., 2016; Krishna et al., 2017; Gupta and Malik, 2015; Dai et al., 2017) and scene graphs (Zellers et al., 2018; Xu et al., 2017).

### 2.5.1 Selecting relationship triplets

We explain here how we selected a set of relationship triplets to be annotated. Each triplet has the form of ⟨`class1`, relationship, `class2`⟩, e.g. ⟨`woman`, playing, `guitar`⟩, ⟨`bottle`, on, `table`⟩. The challenge of selecting triplets lies in balancing several requirements: (i) selecting frequent-enough relationships that can be found in real-world images, (ii) generating enough data to be useful for developing future models, and (iii) selecting non-trivial relationships that cannot be inferred from pure co-occurrence (e.g. a `car` and

a `wheel` in the same image are almost always in a 'part-of' relationship).

To meet these requirements, we select pairs of classes co-occurring sufficiently frequently on the train set, and are not connected by trivial relationships, e.g. ⟨man, wears, `shirt`⟩. We also excluded all 'part-of' relationships, e.g. ⟨`window`, part-of, `building`⟩. To make the task more interesting we make sure several relationships can connect the same pair of objects, so the task cannot be solved simply by detecting a pair of objects: the correct relationship between them must be recognized as well. Finally, we make sure relationship triplets are well defined, i.e. if we select the triplet ⟨`class1`, relationship, `class2`⟩, then we do not include triplet ⟨`class2`, relationship, `class1`⟩, which would make it difficult to disambiguate these two triplets in evaluation. In total, we selected 326 candidate triplets After annotation we found that 287 of them have at least one instance in the train split of Open Images (Tab. 10).

Some examples of the selected relationships are: ⟨man, hits, `tennis ball`⟩, ⟨woman, holds, `tennis ball`⟩, ⟨girl, on, `horse`⟩, ⟨boy, plays, `drum`⟩, ⟨dog, inside of, `car`⟩, ⟨girl, interacts with, `cat`⟩, ⟨man, wears, `backpack`⟩, ⟨chair, at, `table`⟩.

We also introduced further attributes in the dataset, which we represent using the 'is' relationship for uniformity, e.g. ⟨chair, is, wooden⟩. In total we consider 5 attributes corresponding to different material properties ('wooden', 'transparent', 'plastic', 'made of textil', 'made of leather') leading to 42 distinct ⟨object, is, attribute⟩ triplets (all of them turned out to have instances in the train split after verification process).

### 2.5.2 Annotation process

Several prior works tried different schemes for annotating visual relationships. (Lu et al., 2016) proposed a controlled protocol, whereas (Krishna et al., 2017) gives the annotators almost complete freedom. On Open Images we have collected extensive bounding box annotations (Sec. 2.4), so we leverage them in the visual relationship annotation process as follows. For each image and relationship triplet ⟨`class1`, relationship, `class2`⟩ we perform the following steps:

1. Select all pairs of object bounding boxes that can potentially be connected in this relationship triplet. As a criterion we require that their classes match those specified in the triplet and that the two boxes overlap after enlarging them by 20% (our relationships assume objects to have physical contact in 3D space and consequently overlap in their 2D projections).
2. Ask human annotators to verify that the two objects are indeed connected by this relationship.

Note that two objects can be connected by several relationships at the same time, since they are not mutually exclusive.

We report the acceptance rates of the triplet verification process in Table 1. These acceptance rates are rather low, which shows that the selected relationships are hard to predict based just on co-occurrence and spatial proximity of objects. The acceptance rate and the total number of final annotations per relationship triplet is detailed in Figure 11. On average the annotators took 2.6 seconds to verify a single candidate triplet.

|                  | at    | holds | under | all   |
|------------------|-------|-------|-------|-------|
| Acceptance rate  | 58.9% | 27.9% | 2.3%  | 28.2% |

**Table 1 Acceptance rates for the relationship annotation process**. Displaying the relationship with the highest acceptance rate ('at'), the one with the median acceptance rate ('holds'), the one with the lowest acceptance rate ('under'), and the acceptance rate across all triplet candidates.

Note that, due to the annotation process, for each pair of positive image-level labels in an image, we annotate all relationships between all objects with those labels. Therefore, we can have multiple instances of the same relationship triplet in the same image, connecting different pairs of objects (e.g. different men on playing different guitars, and even different chairs *at the same table*, Fig. 22). Note, that we excluded `GroupOf` objects from the annotation process.
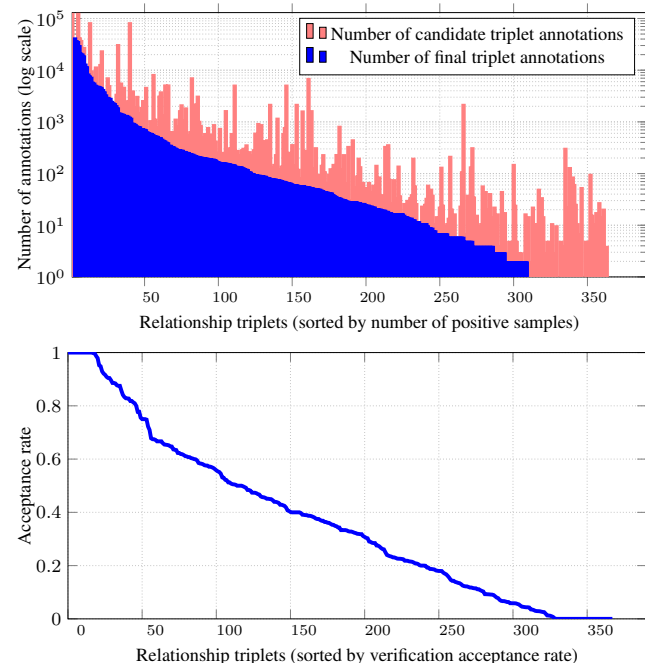


**Fig. 11 Top:** The number of candidate relationship triplet annotations and the number of positively verified ones. The overlap between two object bounding boxes does not guarantee that they are connected by a particular relationship. **Bottom:** Acceptance rate per distinct relationship triplet; note that triplets with acceptance rate 100% have no more than 30 samples in the training split.

## 3 Statistics

The Open Images Dataset consists of 9,178,275 images, split into *train*, *validation*, and *test* (Tab. 2).

|        | Train     | Validation | Test    |
|--------|-----------|------------|---------|
| Images | 9,011,219 | 41,620     | 125,436 |

**Table 2** Split sizes.

As explained in Section 2, the images have been annotated with image-level labels, bounding boxes, and visual relationships, spanning different subsets of the whole dataset. Below we give more detailed statistics about the span, complexity, and sizes of the subsets of images annotated with human-verified image-level labels (Sec. 3.1), with bounding boxes (Sec. 3.2), and with visual relationships (Sec. 3.3).

### 3.1 Human-Verified Image-Level Labels

We assigned labels at the image level for 19,794 classes. Each label can be positive (indicating the class is present in the image) or negative (indicating the class is absent). Figure 4 shows examples and Table 3 provides general statistics.

|                | Train      | Validation | Test      |
|----------------|------------|------------|-----------|
| Images         | 5,655,108  | 41,620     | 125,436   |
| Positive labels| 13,444,569 | 365,772    | 1,105,052 |
| *per image*    | *2.4*      | *8.8*      | *8.8*     |
| Negative labels| 14,449,720 | 185,618    | 562,347   |
| *per image*    | *2.6*      | *4.5*      | *4.5*     |

**Table 3 Human-verified image-level labels**: Split sizes and their label count.

To further study how these labels are distributed, Figure 12 shows the percentage of images with a certain number of positive (solid line) and negative (dashed line) labels. In the train split there are 2.4 positive labels per image on average, while the validation and test splits have 8.8. This discrepancy comes from the fact that we generated candidate labels in the validation and test splits more densely (Sec. 2.3.1) than in the train split (Sec. 2.3.2). Please also note that the distribution of labels for validation and test are the same, since the annotation strategies are the same for both splits.

Some classes are more commonly captured in images than others, and this is also reflected in the counts of annotated labels for different classes. Figure 13 shows the percentage of labels for the top 6,000 classes (sorted by decreasing frequency). As expected, the ~300 most frequent
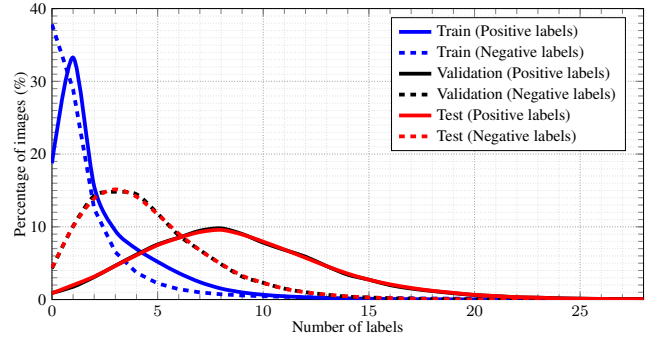


**Fig. 12 Human-verified image-level labels**: Histogram of number of labels per image.

classes cover the majority of the samples for all three splits of the dataset.
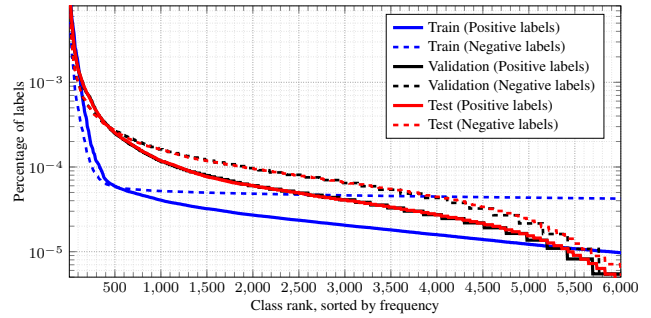


**Fig. 13 Percentage of human-verified image-level labels for each class**. The horizontal axis represents the rank of each class when sorted by frequency, the vertical axis is in logarithmic scale.

As mentioned in Section 2.3.3, label verification is done by annotators from two different pools: internal and crowd-sourced. Table 4 shows the number of human-verified labels coming from each pool. We can see that in train, crowd-sourced labels represent about 20% of all verified labels, whereas for validation and test, they represent less than 1%.

|                          | Train      | Validation | Test      |
|--------------------------|------------|------------|-----------|
| Internal annotators      | 22,351,016 | 547,291    | 1,655,384 |
| Crowd-source annotators  | 5,543,273  | 4,099      | 12,015    |

**Table 4 Internal versus crowd-source human-verified image-level labels**: Number of image-level labels (positive and negative) coming from the two pools of human annotators.

### 3.2 Bounding Boxes

*General statistics*

We annotated bounding boxes around objects of 600 box-

able classes on the whole validation and test splits, and on a subset of the train split (Tab. 5).

|            | Train      | Validation | Test    |
|------------|-----------:|-----------:|--------:|
| Images     | 1,743,042  | 41,620     | 125,436 |
| Boxes      | 14,610,229 | 204,621    | 625,282 |
| *per image*| *8.4*      | *4.9*      | *5.0*   |

**Table 5 Split sizes with annotated bounding boxes**. For each split, number of images and boxes (also normalized per image in italics). These statistics are only over the 600 boxable classes.

Table 6 shows the number of classes, images and bounding boxes in the Open Images Dataset compared to other well-known datasets for object detection: COCO (Lin et al., 2014) (2017 version), PASCAL (Everingham et al., 2015) (2012 version), and ILSVRC (Russakovsky et al., 2015) (2014 detection version). In this comparison we only consider images in Open Images with bounding boxes, not the full dataset. As it is common practice, we ignore the objects marked as difficult in PASCAL. As the table shows, Open Images Dataset is much larger than previous datasets and offers $17\times$ more object bounding boxes than COCO. Moreover, it features complex images with several objects annotated (about the same as COCO on average).

|          | PASCAL | COCO    | ILSVRC-Det All | Dense   | Open Images |
|----------|-------:|--------:|---------------:|--------:|------------:|
| Classes  | 20     | 80      | 200            | 200     | 600         |
| Images   | 11,540 | 123,287 | 476,688        | 80,462  | 1,910,098   |
| Boxes    | 27,450 | 886,284 | 534,309        | 186,463 | 15,440,132  |
| Boxes/im.| 2.4    | 7.2     | 1.1            | 2.3     | 8.1         |

**Table 6 Global size comparison to other datasets**. We take the dataset splits with publicly available ground truth, that is, `train+val` in all cases except Open Images, where we also add the `test` set which is publicly available. Please note that in ILSVRC train, only a subset of ~80,000 images are densely annotated with all 200 classes (~60,000 train and ~20,000 validation). The other images are more sparsely annotated, with mostly one class per image.

To further study how objects are distributed over the images, Figure 14 (left) counts the percentage of images with a certain number of bounding boxes. We can observe that COCO and Open Images are significantly less biased towards single-object images. Figure 14 (right) displays the number of images that contain at least a certain number of bounding boxes. Open Images has significantly more images than the other datasets in the whole the range of number of boxes per image, and especially so at high values, where it covers some regime unexplored before (more than 80 bounding boxes per image, up to 742).

Figure 15 shows some images with a large number of bounding boxes (348, 386, and 743, respectively). In many of these cases, the `GroupOf` attribute could have been used
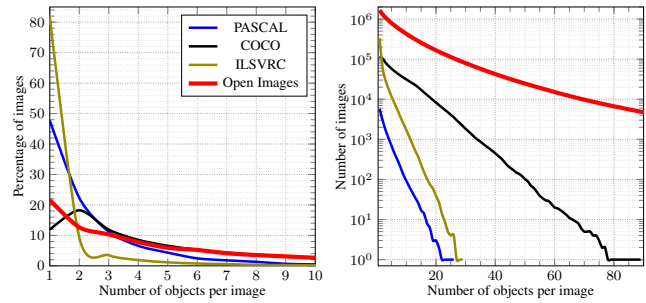


**Fig. 14 Number of objects per image**. Percentage of images with exactly a certain number of objects (left). Number of images with at least a certain number of objects (right). Train set for all datasets.

to reduce the annotation time (e.g. the set of windows marked as `GroupOf` on the right face of the center building). Having some of these extreme cases exhaustively annotated, however, is also useful in practice.

*Box attributes statistics*

As explained in Section 2.4.5, the bounding boxes in Open Images are also labeled with five attributes. Table 7 shows the frequency of these attributes in the annotated bounding boxes.

| Attribute | Occluded | Truncated | GroupOf | Depiction | Inside |
|-----------|---------:|----------:|--------:|----------:|-------:|
| Frequency | 66.06%   | 25.09%    | 5.99%   | 5.45%     | 0.24%  |

**Table 7 Frequency of attributes**: Percentage of boxes with the five different attributes on Open Images train. As reference, the `Crowd` attribute in COCO is present in 1.17% of their boxes.

`Occluded` and `Truncated` are the most common attributes, with a considerable portion of the objects being marked as such. `GroupOf` and `Depiction` are still used in a significant proportion, whereas `Inside` is rare.

Figure 16 displays two images containing boxes labeled with each of the five available attributes. In the left image, the building is viewed from inside, the person is occluded by the stand, and the two busts are depictions of people. The right image shows a group of flowers that is truncated by the picture framing.

*Box class statistics*

Not all object classes are equally common and equally captured in pictures, so the classes in Open Images are not uniformly distributed in their number of instances and through the images. We study both effects below.

Figure 17 (━●━) plots the the number of boxes annotated for each class, sorted by increasing frequency. In order to visually compare to the other datasets with fewer classes, the horizontal axis is shown in logarithmic scale.
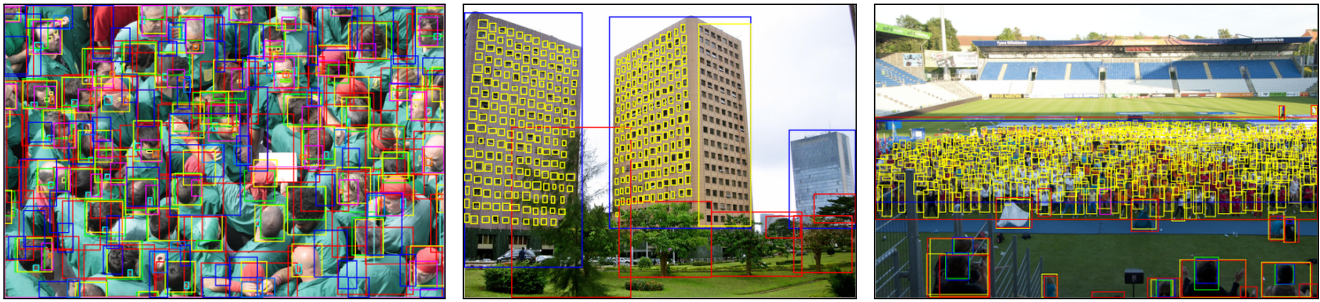
**Fig. 15 Examples of large number of annotated boxes**: Images with 348, 386, and 743, respectively. `GroupOf` could have been used in many of these cases, but nevertheless they still have interest in practice.



**Fig. 16 Examples of box attributes**: `GroupOf`, `Occluded`, `Depiction`, `Truncated`, and `Inside`.
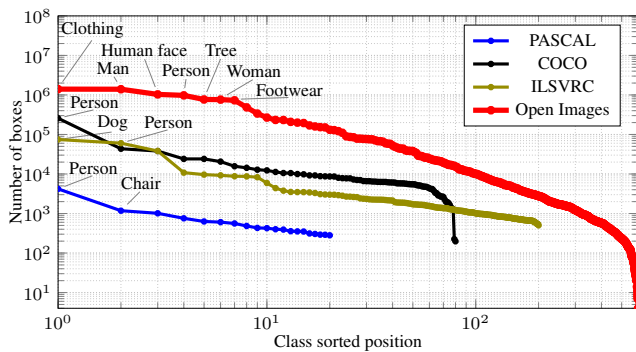


**Fig. 17 Number of boxes per class**. The horizontal axis is the rank of each class when sorted by the number of boxes, represented in logarithmic scale for better readability. We also report the name of the most common classes. Train set for all datasets.

Open Images is generally an order of magnitude larger than the other datasets. There are 11 classes in Open Images with more samples than the largest class in COCO. As a particular example, the `person` class has 257,253 instances in COCO, while Open Images has 3,505,362 instances of the agglomeration of classes referring to person (`person`, `woman`, `man`, `girl`, `boy`)[6].

At the other end of the spectrum, Open Images has 517 classes with more instances than the most infrequent class

---

[6] These are really unique objects: Each object is annotated only with its *leafmost* label, e.g. a `man` has a single box, it is not annotated as `person` also.

in COCO (198 instances), and 417 classes in the case of ILSVRC (502 instances).

Interactions between different object classes are also a reflection of the richness of the visual world. Figure 18 (left) reports the percentage of images with boxes coming from a varying number of distinct classes. We can see that Open Images and COCO have a much richer distribution of images with co-occurring classes compared to ILSVRC and PASCAL, which are more biased to a single class per image.
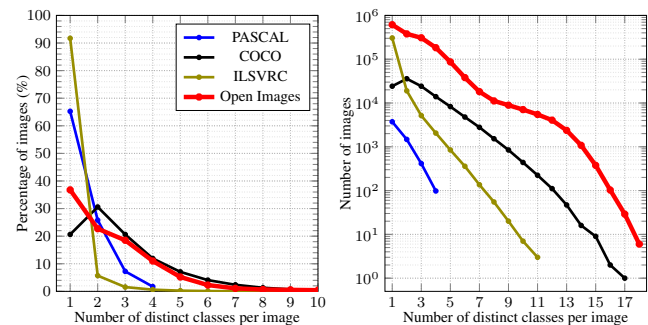


**Fig. 18 Number of distinct classes per image**. Normalized (left) and unnormalized (right) histogram of the number of distinct classes per image. Train set for all datasets.

Figure 18 (right) shows the unnormalized statistics (i.e. with the number of images instead of the percentage). It shows that Open Images has at least one order of magnitude

more images than COCO at any point of the curve. As an example, Open Images has about 1,000 images with 14 distinct classes, while COCO has 20; ILSVRC has no image with more than 11 classes, and PASCAL no more than 4. Figure 19 displays two images with a large number of classes annotated, to illustrate the variety and granularity that this entails.
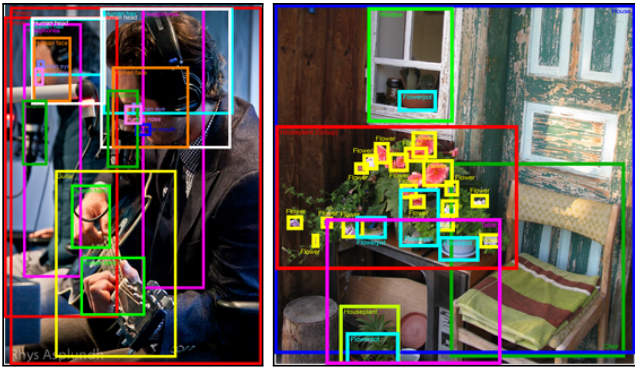


**Fig. 19** Images with a large number of different classes annotated (11 on the left, 7 on the right).

As further analysis, we compute the class co-occurrence matrix in Open Images and sort the pairs of classes in decreasing order. We observe the following patterns. The most co-occurring pairs are human-related classes (`Person`, `Man`, `Woman`) with their parts (`Human face`, `Human arm`, `Human hair`, `Human nose`, etc.) or with accessories (`Clothing`, `Footwear`, `Glasses`, etc.); and other types of objects and their parts (`Car-Wheel`, `House-Window`). Other interesting object pairs co-occurring in more than 100 images are `Drum-Guitar`, `Chair-Desk`, `Table-Drink`, `Person-Book`. Please note that objects co-occurring in an image does not imply them being in any particular visual relationship (analyzed in Sec. 3.3).

*Box size statistics*

Figure 20 displays the cumulative density function of the bounding box sizes in Open Images, PASCAL, COCO, and ILSVRC. The function represents the percentage of bounding boxes (vertical axis) whose area is below a certain percentage of the image area (horizontal axis). As an example, the green lines (——) show that 43% of the bounding boxes in Open Images occupy less than 1% of the image area. Hence, the Open Images Dataset offers a real challenge for object detection, supporting the development and evaluation of future detection models and algorithms.

We compare to two uniform distribution baselines: boxes with uniform area (⋯⋯) or with uniform side length (⋯⋯) (i.e. the square root of their area is uniformly distributed). Interestingly, ILSVRC closely follows the distribution of the
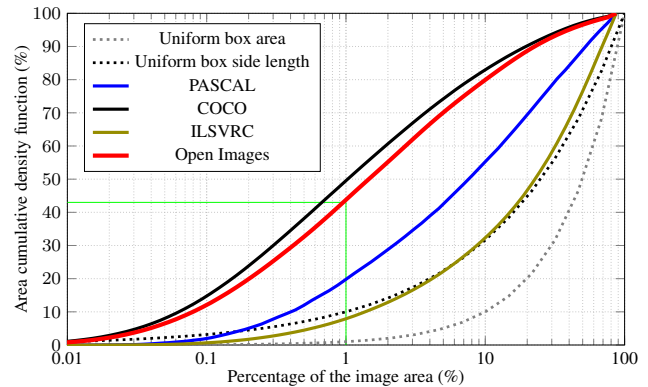


**Fig. 20 Annotated objects area**: Cumulative distribution of the percentage of image area occupied by the annotated objects of PASCAL, COCO, and Open Images; i.e. , percentage of instances whose area is below a certain value. As a baseline, we plot the function corresponding to boxes with uniformly distributed area and side length. We ignore here boxes marked as crowd in COCO and marked as group in Open Images. Train set for all datasets.

latter. In contrast, the other datasets have a greater proportion of smaller objects: Open Images has a similar distribution to COCO, both having many more small objects than PASCAL.

*Box center statistics*

As another way to measure the complexity and diversity of the boxes, Figure 21 shows the distributions of object centers[7] in normalized image coordinates for Open Images and other related datasets. The Open Images train set, which contains most of the data, shows a rich and diverse distribution of a complexity in a similar ballpark to that of COCO. Instead, PASCAL and ILSVRC exhibit a simpler, more centered distribution. This confirms what we observed when considering the number of objects per image (Fig. 14) and their area distribution (Fig. 20).

*Validation and test V5*

The number of boxes per image (Tab. 5) in Open Images V4 is significantly higher in the train split than in validation and test. In the next version of Open Images (V5) we increased the density of boxes for validation and test to be closer to that of train (Tab. 8).

|  | Train | Validation | Test |
|---|---|---|---|
| Boxes V4 | 14,610,229 | 204,621 | 625,282 |
| *per image* | *8.4* | *4.9* | *5.0* |
| Boxes V5 | 14,610,229 | 303,980 | 937,327 |
| *per image* | *8.4* | *7.3* | *7.5* |

**Table 8 Number of boxes for Open Images Dataset V4 versus V5**.

---

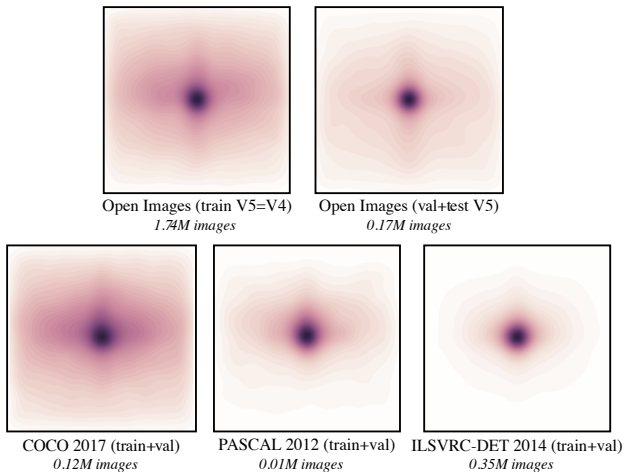[7] We thank Ross Girshick for suggesting this type of visualization.

**Fig. 21 Distribution of object centers** for various splits of Open Images and other related datasets.

Fig. 21 shows the distributions of object centers. While the smaller val and test splits are still simpler than train, they are considerably richer than ILSVRC and also slightly better than PASCAL.

### 3.3 Visual Relationships

The Open Images Dataset was annotated in a very controlled manner. First, we produced image-level labels verified by human annotators. Afterwards, we annotated bounding boxes on *all instances* of each positive image-level label for 600 classes. Now we expanded the annotations beyond object bounding boxes: we precisely defined a set of relationships between objects and then verified their presence for each pair of annotated objects in the dataset (Sec. 2.5). In the end of the process we obtained 375k annotations of 329 distinct relationship triplets, involving 57 different object classes. Figure 22 shows example annotations and Table 9 contains datasets statistics per split (train, validation, test).

|  | Train | Val | Test |
|---|---|---|---|
| Number of VRD annotations | 374, 768 | 3, 991 | 12, 314 |

**Table 9** Number of annotated visual relationship instances for the train, validation and test splits of Open Images.[9]

On the other end of the annotation spectrum is data collection as proposed by the creators of Visual Genome (VG) and Visual Relationship Detection (VRD) datasets (Krishna et al., 2017; Lu et al., 2016). Their focus was on obtaining

as much variety of relationships as possible by asking annotators to give a free-form region description, and annotate objects and relationships based on those descriptions. The annotations from several annotators were then merged and combined using various language and quality models.

The difference in the two approaches naturally leads to difference in the properties of the two datasets: while VG and VRD contain higher variety of relationship prepositions and object classes (Tab. 10) they also have some shortcomings. First, previous work shows that many of those are rather obvious, i.e. ⟨ window, on building ⟩ (Zellers et al., 2018). Table 11 compares the top-10 most frequent relationship triplets in all three datasets: in both VG and VRD the most frequent relationships can be predicted from object co-occurrence and spatial proximity, while Open Images is more challenging in this respect. Second, as follows from the free-form annotation process and lack of precise predefinitions, the annotations on VG and VRD contain multiple relationships with the same semantic meaning: for example, the difference between relationships 'near' and 'next to' is not clear. This leads to annotation noise as multiple instances of conceptually the same relationship have different labels. Since Open Images annotations were collected in a very controlled setting this kind of noise is much lower. Finally, in VG and VRD annotations within an image are sometimes incomplete (e.g. if there are two chairs at a table in the same image, only one of them might be annotated). Instead, thanks to the controlled annotation process for image-level labels, boxes, and relationships, in Open Images for each image it is possible to know exactly if two objects are connected by a certain relationship or not. This makes Open Images better suited for evaluating the performance of visual relationship detection models, and also facilitates negative samples mining during training.

As one can expect from object class distribution on Open Images, the distribution of the number of relationship annotations among triplets is highly imbalanced (Fig. 23). Hence, the Open Images dataset includes both rare and very frequent relationship triplets. This suggests that to be able to effectively detect triplets that have very small number of annotations, it will not be enough to train a monolithic detector for each triplet. We expect that a successful detector will have to be compositional.

Figures 24 and 25 provide a comparison between the number of annotations in Open Images vs in VG/VRD for the semantic triplets they have in common (considering only two-object relationships, not attributes). As the plots shows, the Open Images Dataset has more annotations for several triplets than VG/VRD, which shows it can complement them[10]. Moreover, Open Images contains new relationship triplets
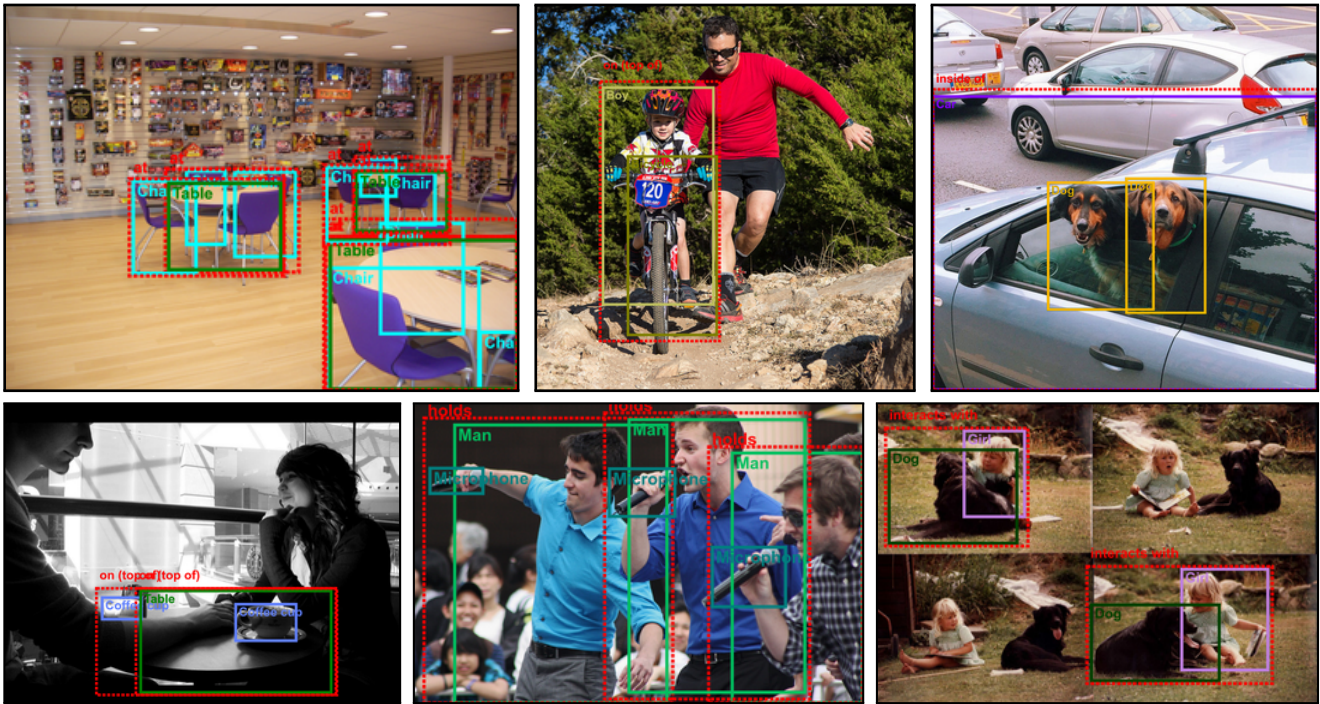
**Fig. 22 Examples of positively verified relationships**. Note how we annotated all instances of a relationship triplet (e.g. multiple ⟨ man, holds, microphone ⟩ in the same image).
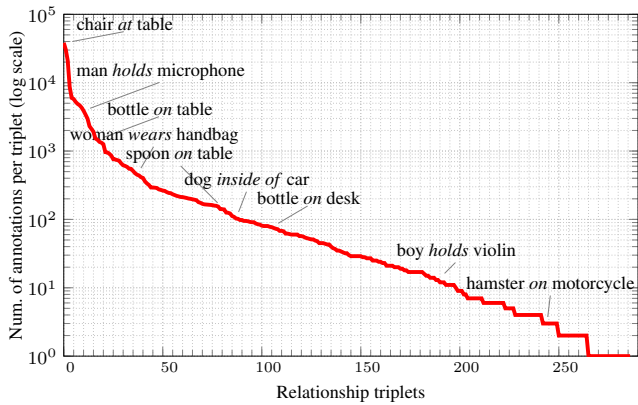


**Fig. 23** Number of annotations per triplet on Open Images (two-object relationships only, without attributes).
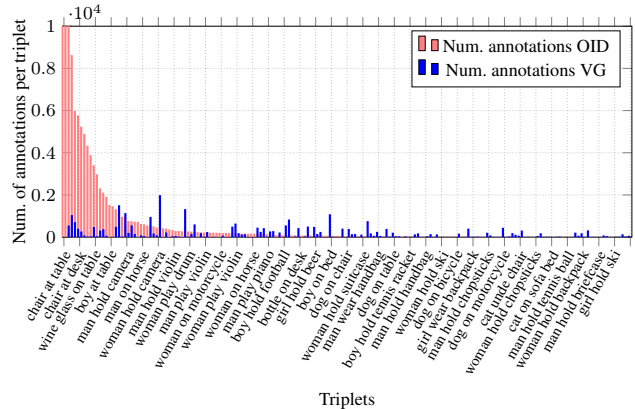


**Fig. 24 Comparison of the number of triplet annotations** on Open Images dataset vs. Visual Genome dataset for 195 triplets found in common (two-object relationships only, without attributes).

than are not in VG at all, e.g. ⟨ man, play, flute ⟩, ⟨ dog, inside of, car ⟩, ⟨ woman, holds, rugby ball ⟩.

In summary, Open Images visual relationship annotations are not as diverse as in VG and VRD, but are better defined, avoid obvious relationships, have less annotation noise, and are more completely annotated. Moreover, Open Images offers some complementary annotations to VG/VRD, both by the number of samples for some the triplets they have in common, and by some entirely new triplets.

---

inconsistent relationship names, we use loose string matching to match relationships

## 4 Quality

### 4.1 Quality of bounding boxes

We performed an extensive analysis of the quality of bounding box annotations. We had a human expert examine 100 images for each of the first 150 boxable classes sorted by alphabetical order, containing a total of more than 26,000 boxes. We measured the quality of both the boxes and their attributes.

Results for box quality are shown in Table 12. Both precision and recall are very high at 97.7% and 98.2%, re-

| | Num. classes (/ num. attributes) | Num. distinct triplets | Num. annotations |
|---|---|---|---|
| Visual Relationship Detection (Lu et al., 2016) | 100 | 6,672 | 30,355 |
| Visual Genome (Krishna et al., 2017) | 67,123 / 4,279 | 727,635 | 2,578,118 |
| *two-object relationships* | 65,398 | 675,378 | 2,316,104 |
| *attributes* | 7,100 / 4,279 | 52,257 | 262,014 |
| Open Images | 57 / 5 | 329 | 374,768 |
| *two-object relationships* | 57 | 287 | 180,626 |
| *attributes* | 23 / 5 | 42 | 194,142 |

**Table 10** Comparison with the existing visual relationship detection datasets.
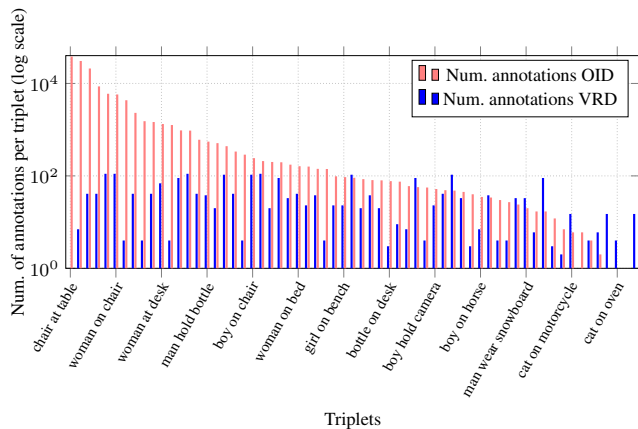


**Fig. 25 Comparison of the number of triplet annotations** on the Open Images Dataset versus the VRD dataset for 62 triplets found in common (two-object relationships only, without attributes).

| VRD dataset | Visual Genome | Open Images |
|---|---|---|
| person *wear* shirt | window *on* building | chair *at* table |
| person *wear* pants | clouds *in* sky | man *at* table |
| person *next to* person | man *wearing* shirt | woman *at* table |
| person *wear* jacket | cloud *in* sky | man *on* chair |
| person *wear* hat | sign *on* pole | woman *on* chair |
| person *wear* glasses | man *wearing* hat | chair *at* desk |
| person *has* shirt | leaves *on* tree | man *holds* guitar |
| person *behind* person | man *wearing* pants | man *plays* guitar |
| person *wear* shoes | man *has* hair | chair *at* coffee table |
| shirt *on* person | building *has* window | girl *at* table |

**Table 11** Top-10 most frequent relationships in the VRD, Visual Genome and Open Images datasets. Note that some of the Open Images relations are not mutually exclusive (e.g. "man holds guitar" and "man plays guitar"). In these cases, we have annotated all relationships that occur in each particular sample (see Section 2.5.2).

spectively. The missing precision is mostly caused by boxes which are geometrically imprecise (1.1%), and boxes with wrong semantic class labels (1.1%). Imprecise boxes are quite evenly spread over classes. However, while half of the classes have fewer than 1% semantic errors, other (often rare) classes have more (Fig. 26). Most notably, the three most problematic classes are `bidet` (86% errors, confused with toilets), `cello` (55% errors, confused with violins), and `coffee table` (35% errors, confused with other kinds

of tables). The first two mistakes are caused by cultural differences. `coffee table` is an ambiguous class.

| precision | | | recall |
|---|---|---|---|
| | 97.7% | | 98.2% |
| imprecise | wrong class | multiple objects | |
| 1.1% | 1.1% | 0.1% | |

**Table 12 Analysis of box quality**. Conditioned on a given class label for an image, we report precision and recall. We break down precision errors into three different types: an geometrically imprecise box, a box with the wrong class label, and a box which unjustifiably captures multiple objects.



**Fig. 26 Percentage of boxes which have a semantic error** for each of the 150 examined classes. Every fifth class name is displayed on the horizontal axis. Moreover, we provide the names of the 9 classes with the highest percentage of errors directly on the curve.

Table 13 shows results for attribute quality. Precision and recall are in the high nineties for most attributes. Especially the `Occluded` and `Truncated` attributes are very accurately annotated. For `Depiction`, recall is 92%. For `Inside`, precision is 67%, which is mostly caused by several `Bell peppers` incorrectly having this attribute when it was inside a container such as a shopping cart. However, the `Inside` attribute is extremely rare (0.4% of all boxes)

and only relevant for a few classes. The most frequent such class is `Building`, for which precision and recall are good (100% and 83%, respectively).

| Attribute | Precision | Recall |
|---|---|---|
| GroupOf | 94.2% | 95.3% |
| Occluded | 98.6% | 98.4% |
| Truncated | 99.7% | 97.0% |
| Depiction | 96.7% | 92.2% |
| Inside | 66.7% | 90.3% |

**Table 13 Precision and recall for attributes**, with the error rates specified per attribute type.

## 4.2 Geometric agreement of bounding boxes

Another way to measure the quality of bounding boxes is to draw them twice by different annotators, and then measuring their geometric agreement. We did that for both Extreme Clicking (Sec. 2.4.2) and Box Verification Series (Sec. 2.4.3).

*Extreme clicking.* We randomly selected 50,000 boxes produced with extreme clicking, and had annotators redraw the box (without seeing the original box). We then measured human agreement as the average intersection-over-union (IoU) between the original box and the redrawn box on the same object. We found this to be 0.87, which is very close to the human agreement of 0.88 IoU on PASCAL (Everingham et al., 2015) reported by (Papadopoulos et al., 2017). The slight difference is mainly caused by objects being generally smaller in Open Images (Fig. 20).

*Box Verification Series.* For the boxes produced with box verification series, we had 1% re-annotated using extreme clicking (again without showing the original box). We then measured IoU between the original boxes and the newly manually drawn boxes. We found this to be 0.77 on average. As expected, this is higher than the threshold of IoU$> 0.7$ for which we trained the annotators, and lower than the extreme clicking agreement of 0.87. We underline that 0.77 is widely considered as a good geometric accuracy (e.g. the COCO Challenge calls IoU $> 0.75$ a "strict" evaluation criterion). To give a better feeling of the average quality of these boxes, Fig. 27 shows two examples where a drawn box and box produced by box verification series have IoU$=$ 0.77.

*Detectors: extreme clicking vs box verification series.* The extreme clicking boxes are more accurate than those produced by box verification series. But how does this influence contemporary object detection models? To answer this question, we make use of the 1% of re-annotated data from box



**Fig. 27 Two examples of matching boxes with** $IoU = 0.77$, the average agreement between verified and drawn boxes. Verified boxes are in green, drawn boxes are in red. As these example shows, the verified boxes cover the object very well, but not perfectly.

verification series. This means that for the exact same object instances, we have both verified boxes and manually drawn boxes. We train a Faster-RCNN (Ren et al., 2015) model based on Inception-ResNet-V2 (Szegedy et al., 2017) on each set and measure performance on the Open Images test split. We measure performance using the Open Images Challenge metric $mAP_{OI}$ (Sec. 5.2.1), which is a modified version of mean Average Precision commonly used for object detection (Everingham et al., 2010). Interestingly, we found that the difference in detection performance was smaller than 0.001 $mAP_{OI}$. We conclude that boxes produced by box verification series make perfectly useful training data for contemporary detectors, as they lead to the same performance as training from manually drawn boxes.

## 4.3 Recall of boxable image-level labels

As described in Sec. 2.3, we obtained image-level labels by verifying candidate labels produced automatically by a classifier. Here we estimate the general recall of this process for the 600 boxable classes. We randomly sampled 100 images and inspected each image independently by two human experts to identify all instances of boxable classes which were not annotated. This was done by displaying each image with all existing box annotations overlaid. For each non-boxed object, the expert typed multiple free-form words, each of which was mapped to the closest five Open Images boxable classes through Word2Vec (Mikolov et al., 2013). Based on these, the expert then decided whether the object indeed belonged to a boxable class, and recorded it as a missing object. Additionally, the object could be marked as 'difficult' according to the PASCAL standards (Everingham et al., 2010), i.e. very small, severely occluded or severely truncated. Afterwards, we took the set union of all labels of all missing objects recorded by the two experts. We removed existing image-level labels from this set, to cover for the rare case when an object instance was not boxed even though its image-level label was available (high recall in Tab. 12). This results in the final set of classes present in the image but for which we do not have an image-level label.

When considering really all objects, the recall of image-level labels is 43% in Open Images. When disregarding 'difficult' objects, the recall is 59%. While this is lower than the estimated recall of 83% reported for COCO (Lin et al., 2014), Open Images contains $7.5\times$ more classes, making it much harder to annotate completely. Importantly, this lack of complete annotation is partially compensated by having explicit negative image labels. These enable proper training of discriminative models. Finally, we stress that for each positive image-level label we annotated bounding boxes for *all* instances of that class in the image. Along that dimension, the dataset is fully annotated (98.2% recall, Tab. 12).

## 5 Performance of baseline models

In this section we provide experiments to quantify the performance of baseline models for image classification (Sec. 5.1) and object detection (Sec. 5.2) on the Open Images Dataset.

### 5.1 Image Classification

Image classification has fostered some of the most relevant advances in computer vision in the last decade, bringing new techniques whose impact has reached well beyond the task itself (Krizhevsky et al., 2012; Szegedy et al., 2015; Ioffe and Szegedy, 2015; He et al., 2016; Szegedy et al., 2017). Here we train an image classification model with Inception-ResnetV2 (Szegedy et al., 2017), a widely used high capacity network, and we empirically measure:

- The impact of the number of human-verified labels on the quality of a classifier.
- The impact of using negative human-verified labels.
- Classification performance when restricted to boxable classes only.

*Training*

For our experiments we use the model described in Section 2.3.1 to produce machine-generated labels on the train split. We consider each label predicted with confidence score above $0.5$ as positive. This way we have two sets of labels for the train split: these machine-generated labels, and the human-verified labels as discussed in Section 2.3.3.

In these experiments we consider only classes with at least 100 positive human-verified examples in the train split (7,186 classes).

We first pre-train the network from scratch using the machine-generated labels. We then fine-tune it with a mix of 90% human-verified labels and 10% machine-generated labels.

*Evaluation*

For each class we calculate Average Precision on the test set, for the 4,728 classes that are both in the trainable label set and have $> 0$ samples in the test set. During evaluation we take into account that the ground-truth annotations are not exhaustive (Sec. 2.3), and do not penalize the model for predicting classes for which we do not have human verification on the test set. This metric is discussed in detail for Open Images in (Veit et al., 2017).

*Number of human-verified labels*

To measure the impact of the human-verified labels we repeat the fine-tuning stage described above to build classification models using a varying fraction of the human-verified labels. In all case we start from the model pre-trained on machine-generated labels on the entire train split. We then fine-tune on human-verified labels from random subsets of the train split containing 1%, 10%, 25%, 50%, 75%, and 100% of all images.

As shown in Fig. 28, mAP increases as we increase the amount of human-verified labels, demonstrating that they directly improve model performance. The absolute number of human-verified labels can be calculated using the values in Table 3.

*Impact of negative labels*

To measure the impact of negative image-level labels we repeat the above experiment but train from positive labels only (Fig. 28). For this we train our classification model while ignoring human-verified negative labels in the loss. Instead, we use as implicit negatives all labels that are not human-verified as positive (including candidate labels generated by the image classifier, and all other labels that are missed by it, Sec. 2.3). We observe that, as saturation starts to occur when using a large number of positive labels, negative labels start to improve the model significantly. Hence, this experiment shows the value of explicit, human-verified negative labels.

*Boxable classes*

We also report mAP for the 600 boxable classes (using the same models as before). As shown in Fig. 28, mAP for these boxable classes is generally higher than mAP for all classes. Boxable classes are generally easier for classification tasks, as they are all concrete objects defined clearly by their visual properties (as opposed to some of the classes in the wider 19,794 set, e.g. love and birthday). Also, they are mostly basic-level categories, whereas the wider set contains many fine-grained classes which are harder to distinguish (e.g. breeds of dogs). As before, we observe that a larger number of human-verified labels translates to a higher value of mAP.
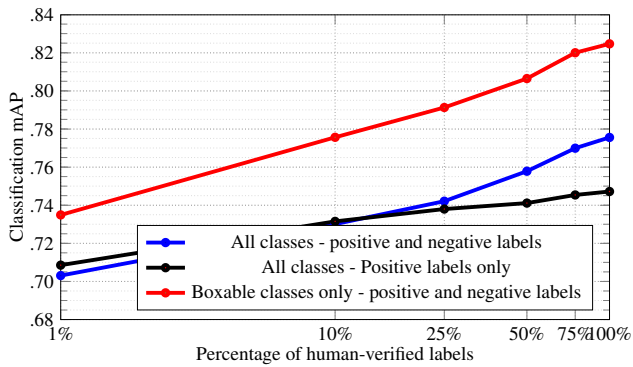
**Fig. 28 Classifier performance versus amount of human-verified labels** in terms of the percentage of all available labels.

## 5.2 Object detection

The advent of object detection came in the form classifiers applied densely to windows sliding over the image (e.g. based on boosting (Viola and Jones, 2001a,b) or Deformable Part Models (Felzenszwalb et al., 2010a,b)). To reduce the search space, the concept of "object proposals" (Alexe et al., 2010, 2012; Uijlings et al., 2013) was then introduced, which enable to work on just a few thousand windows instead of a dense grid.

R-CNN (Girshick et al., 2014) brought the advances in image classification using deep learning to object detection using a two-stage approach: classify object proposal boxes (Uijlings et al., 2013) into any of the classes of interest. This approach evolved into Fast R-CNN (Girshick, 2015) and Faster R-CNN (Ren et al., 2015), which generates the proposals with a deep network too. Faster R-CNN stills provide very competitive results today in terms of accuracy.

More recently, single-shot detectors were presented to bypass the computational bottleneck of object proposals by regressing object locations directly from a predefined set of anchor boxes (e.g. SSD (Liu et al., 2016) and YOLO (Redmon et al., 2016; Redmon and Farhadi, 2017)). This typically results in simpler models that are easier to train end-to-end.

In this section we evaluate the performance of two modern object detectors on Open Images (the two-stage Faster-RCNN (Ren et al., 2015) and the single-shot SSD (Liu et al., 2016)). We start by defining an evaluation metric that takes into account the characteristics of Open Images in Sec. 5.2.1. Then we detail our evaluation setup and report results exploring various model architectures and training set sizes in Section 5.2.2.

### 5.2.1 Evaluation metric

The standard metric used for object detection evaluation is PASCAL VOC 2012 mean average precision (mAP) (Everingham et al., 2012). However, this metric does not take into account several important aspects of the Open Images Dataset: non-exhaustive image-level labeling, presence of class hierarchy and group-of boxes. We therefore propose several modifications to PASCAL VOC 2012 mAP, which are discussed in detail below.

*Non-exhaustive image-level labeling.* Each image is annotated with a this is not the spirit of the point we are making, it is a boundary case; and anyway sets can be empty! set of positive image-level labels, indicating certain classes are present, and negative labels, indicating certain classes are absent. All other classes are unannotated. Further, for each positive image-level label, every instance of that object class is annotated with a ground-truth bounding box. For fair evaluation we ignore all detections of unannotated classes. A detection of a class with a negative label is counted as false positive. A detection of a class with a positive label is evaluated as true positive or false positive depending on its overlap with the ground-truth bounding boxes (as in PASCAL VOC 2012). For a detection to be evaluated as true positive, its intersection-over-union with a ground-truth bounding box should be greater than $0.5$.

*Class hierarchy.* Open Images bounding box annotations are created according to a hierarchy of classes (Section 2.4.4). For a leaf class in the hierarchy, AP is computed as normally in PASCAL VOC 2012 (e.g. 'Football Helmet'). In order to be consistent with the meaning of a non-leaf class, its $AP$ is computed involving all its ground-truth object instances and all instances of its subclasses. For example, the class `Helmet` has two subclasses (`Football Helmet` and `Bicycle Helmet`). These subclasses in fact also belong to `Helmet`. Hence, $AP_{\texttt{Helmet}}$ is computed by considering that the total set of positive `Helmet` instances are the union of all objects annotated as `Helmet`, `Football Helmet`, and `Bicycle Helmet` in the ground-truth. As a consequence, an object detection model should to produce a detection for each of the relevant classes, even if each detection corresponds to the same object instance. For example, if there is an instance of `Football Helmet` in an image, the model need to output detections for both `Football Helmet` and for `Helmet` in order to reach 100% recall (see the semantic hierarchy visualization in Fig.5). If only a detection with `Football Helmet` is produced, one true positive is scored for `Football Helmet` but the `Helmet` instance will not be detected (false negative).

*Group-of boxes.* A group-of box is a single box containing several object instances in a group (i.e. more than 5 instances which are occluding each other and are physically touching). The exact location of a single object inside the group is unknown.

We explore two ways to handle group-of boxes. These can be explained in a unified manner, differing in the value of a parameter weight $w \in \{0, 1\}$. If at least one detection is inside a group-of box, then a single true positive with weight $w$ is scored. Otherwise, the group-of box is counted as a single false negative with the same weight $w$. A detection is inside a group-of box if the area of intersection of the detection and the box divided by the area of the detection is greater than 0.5. Multiple correct detections inside the same group-of box still count as a single true positive.

When $w = 0$ group-of boxes act like ignore regions: the detector is not required to detect them, and if it does output detections inside them, they are ignored. Instead, when $w = 1$ the detector is required to output at least one detection inside a group-of box. In our final evaluation metric, we use $w = 1$.

*Effects.* To evaluate the effect of proposed customized evaluation metric we show results on the Faster-RCNN detector (Ren et al., 2015) with Inception-ResNetV2 backbone (Szegedy et al., 2017) using various versions of the metric. The details of training are given in the following Section 5.2.2. As Figure 29 shows, the biggest difference is caused by ignoring detections of unannotated classes, and thus taking into account the non-exhaustiveness of the annotations. Without this, correct detections of objects from unannotated classes would be wrongly counted as false negatives.
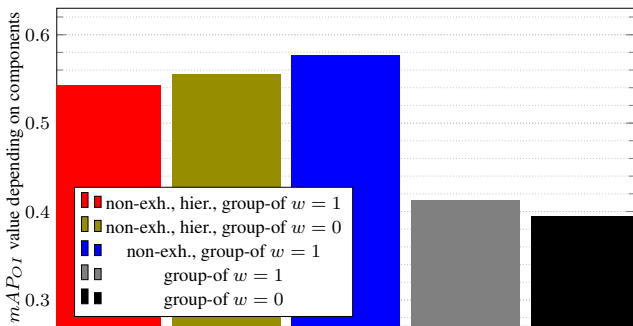
**Fig. 29 Effect of the components of the Open Images metric**. The full metric (non-exh., hier., group-of $w = 1$) and the effect of its components: non-exhaustive labeling, presence of hierarchy, group-of box weight $w$. The black bar represents a metric very close to the standard PASCAL VOC 2012 mAP (with the only addition of ignoring detections inside a group-of box).

### 5.2.2 Performance evaluation

We evaluate two modern object detection models with different capacities. The first model is Faster-RCNN (Ren et al., 2015) with an Inception-ResNetV2 backbone, which performs feature extraction (Szegedy et al., 2017). The second

model is SSD (Liu et al., 2016) with MobileNetV2 (Sandler et al., 2018) feature extractor with depth multiplier 1.0 and input image size $300 \times 300$ pixels. We report in Table 14 the number of parameters and inference speed for each detection model.

| Detector | Number of parameters | Inference time (s) |
|---|---|---|
| Faster-RCNN with Inception-ResNetV2 | 63.947.366 | 0.455 |
| SSD with MobileNetV2, dm=1.0 | 14.439.167 | 0.024 |

**Table 14 Detector capacity**: number of parameters and inference speed measured on a Titan X Pascal GPU.

We consider four increasingly large subsets of the Open Images train set, containing 10k, 100k, 1M and 14.6M bounding boxes. We train both detectors on exactly the same subsets and test on the publicly released Open Images test set. All feature extractors are pre-trained on ILSVRC-2012 (Russakovsky et al., 2015) for image classification until convergence. Then, the models are trained for object detection on Open Images for 8M-20M steps until convergence on 8-24 NVidia GPUs (Tesla P100, Tesla V100). For the Faster-RCNN architecture we use momentum optimizer (Qian, 1999), while for the SSD architecture we used RMS-prop All hyperparameters are kept fixed across all training sets.

Figure 30 reports the results for each combination of deteciton model and training subset size. Generally, the performance of all detectors continuously improves as we add more training data. Faster-RCNN with InceptionV2 improve all the way to using all 14.6M boxes, showing that the very large amount of training data offered by Open Images is indeed very useful. The smaller SSD with MobileNetV2 detector saturates at 1M training boxes. This suggests that for smaller models Open Images provides more than enough training data to reach their performance limits.
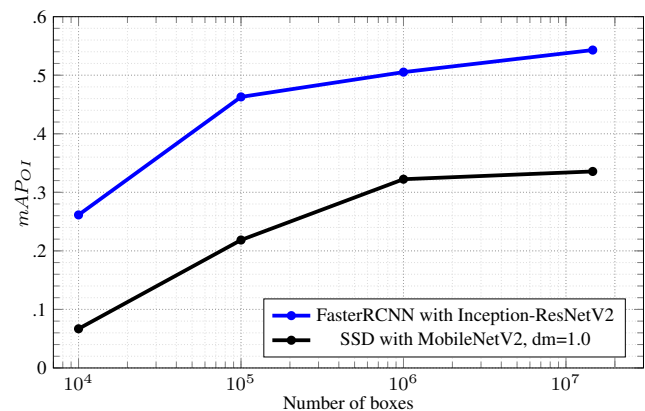
**Fig. 30 Detector performance vs training set size**.

## 5.3 Visual relationship detection

Many works have been proposed to tackle the visual relationship detection task in a fully supervised scenario (Lu et al., 2016; Liang et al., 2017; Dai et al., 2017; Zhang et al., 2017a; Li et al., 2017), in a weakly supervised setting (Peyre et al., 2017; Zhang et al., 2017b), or focusing on human-object relations (Gupta et al., 2009; Yao and Fei-Fei, 2010; Prest et al., 2012). Recently, high-performing models based on deep convolutional neural networks are dominating the field (Gupta and Malik, 2015; Gkioxari et al., 2018; Gao et al., 2018; Kolesnikov et al., 2018).

In this section we evaluate two frequency baselines (Zellers et al., 2018) as well as a state-of-the-art visual relationship detection model (Kolesnikov et al., 2018).

*Tasks and evaluation.* Traditionally, three main tasks were considered in the visual relationship detection community: relationship detection, phrase detection and preposition detection (Lu et al., 2016). The first task is the most challenging one, while the other tasks are a relaxation of it. The performance of visual relationship detection models is often measured either as Recall@50, @100, etc (Lu et al., 2016), or using mean average precision (mAP) as for object detection (Gupta and Malik, 2015). However, until recently VRD datasets did not provide exhaustive annotations (except human-centric datasets (Gupta and Malik, 2015)). Unfortunately this makes the mAP metric deliver an overly pessimistic assessment, since correct model predictions are sometimes scored as false positives due to missing annotations. Open Images provides image-level labels that indicate if a given object class was annotated in an image (either as present or absent, Sec. 2.3). If a class was annotated, then all its object instances are annotated with bounding boxes, and also all its occurrences in visual relationship triplets are also exhaustively annotated. Hence, mAP can be computed in similar manner as object detection mAP, while ignoring predictions on images where annotation are not present according to the image-level labels (Section 5.2.1). Thus, there is no risk of incorrectly over-counting false-positives.

In the next paragraphs we provide performance evaluation of several baselines using two metrics:

– mAP for the visual relationship detection task only (taking into account image-level labels to not penalize correct predictions if ground-truth annotations are missing).
– the Open Images Challenge metric[11], which is a weighted average of three metrics: mAP for relationship detection, mAP for phrase detection, and Recall@50.

*Frequency baselines.* We compute two frequency baselines inspired by (Zellers et al., 2018). As in (Zellers et al., 2018),

---

[11] https://storage.googleapis.com/openimages/web/evaluation.html

we name them FREQ and FREQ-OVERLAP. Let $S$ indicate the subject, $O$ the object and $P$ the relationship preposition connecting two objects. We first model the probability distribution $p(S, P, O|I)$ that a triplet $\langle S, P, O \rangle$ is a correct visual relationship in the input image $I$, using the chain rule of probability. The joint probability distribution can be decomposed into:

$$p(S, P, O|I) = p(S|I) \cdot p(P|O, S, I) \cdot p(O|S, I) \qquad (2)$$

In the simplest case $p(P|O, S, I)$ can be computed from the training set distribution as the prior probability to have a certain relationship given a bounding box from a subject $S$ and object $O$, without looking at the image, i.e. $p(P|O, S, I) = p(P|O, S)$. For the FREQ baseline it is computed using all pairs of boxes in the train set. FREQ-OVERLAP instead is computed using only overlapping pairs of boxes. Further, assuming the presence of $O$ is independent of $S$, then $p(O|S, I) = p(O|I)$.

To compute the $p(O|I)$ and $p(S|I)$ factors we use the FasterRCNN with Inception-ResNetV2 object detection model from Section 5.2, and RetinaNet (Lin et al., 2017) with ResNet50 that serves as a base model for BAR-CNN. After the set of detections is produced, we derive the final score for each pair of detections according to Eq. (2) and using the prior (FREQ baseline). For the FREQ-OVERLAP baseline, only overlapping pairs of boxes are scored using the corresponding prior. In a summary, these baselines use an actual object detector applied to the input image to determine the location of the subject and object boxes, but then determines their relationship based purely on prior probabilities, as learned on the training set.

*BAR-CNN baseline.* The BAR-CNN model (Kolesnikov et al., 2018) is a conceptually simple model that first predicts all potential subjects in an image and then uses an attention mechanism to attend to each subject in turn and predict all objects connected with it by a relationship. This model is shown to deliver state-of-the-art results despite its simplicity. We train BAR-CNN with ResNet50 backbone and focal loss (Lin et al., 2017) on the training set with bounding boxes of and then fine-tune it for the visual relationship detection task using visual relationship annotations. In contrast to the frequency baselines, BAR-CNN considers the input image also for predicting the relationship, and detects the object conditioned on the subject.

### 5.3.1 Performance evaluation

The evaluation results are presented in Table 15. The FREQ and FREQ-OVERLAP baselines score relatively low on the task, even when based on a strong object detector. This indicates that visual relationship detection requires more than simply object detection plus relationship priors. BAR-CNN

instead performs much better than the frequency baselines. That indicates that there is a lot of extra visual information needed to correctly identify visual relationships on an image. The result of BAR-CNN can be regarded as a reference for further improvements on the Open Images dataset.

| Baseline | mAP | score |
|---|---|---|
| FREQ | 5.36 | 18.93 |
| FREQ-OVERLAP | 8.19 | 21.47 |
| FREQ (RetinaNet+ResNet50) | 5.68 | 19.01 |
| FREQ-OVERLAP (RetinaNet+ResNet50) | 8.12 | 21.02 |
| BAR-CNN (RetinaNet+ResNet50) | 14.63 | 27.60 |

**Table 15 VRD baselines**: performance of the VRD models on the test set of Open Images.

## 6 The Power of a Unified Dataset

Unification is one of the distinguishing factors of the Open Images dataset, in that the annotations for image classification, object detection, and visual relationship detection all coexist on the same images. In this section we present two experiments that take advantage of the different types of annotations present in the same images.

### 6.1 Fine-grained object detection by combining image-level labels and object bounding boxes

The Open Images dataset contains 19,794 classes annotated at the image-level and 600 classes annotated at the box-level. Bounding boxes provide more precise spatial localization but image-level labels are often semantically more fine-grained and specific. Since all classes in Open Images are part of a unified semantic hierarchy, we can find the image-level classes that are more specific (children) than a certain bounding-box class (parent). As an example, there are image-level classes such as Volkswagen or Labrador that are more specific than the bounding-box classes Car or Dog, respectively. The experiment in this section shows how we can create a fine-grained object detector (e.g. Volkswagen or Labrador) by combining the two types of annotations.

*Creating fine-grained detection data.* Given a bounding-box class cls, we denote the set of all image-level classes more specific than cls as C(cls). We then look for images where there are boxes of class cls and image-level labels of any class of C(cls). In those images which have only one bounding box of class cls, we transfer the more fine-grained labels C(cls) to it. This transfer is safe, as there is only one possible object in that image.

We looked for bounding box classes with a significant number of more specific image-level classes and selected the following four to experiment with: Car, Flower, Cat, and Dog. We use the procedure above to create fine-grained box labels for these four classes. Statistics are presented in Table 16. Finally, we use stratified sampling to divide our data into 90% training images and 10% test images.

*Experimental setup.* We evaluate on fine-grained classes which have at least 4 training samples and at least 1 test sample. We train a single Faster-RCNN detector (Ren et al., 2015) with an Inception-ResNetV2 backbone (Szegedy et al., 2017) (like in Sec. 5.2.2) on the fine-grained classes. We apply this detector on the test set and report the Average Precision at IoU > 0.5, averaged over the fine-grained classes within each general class (mAP over Car, Cat, Dog, Flower).

We also report the performance of three baselines, all based on the same Faster-RCNN architecture as above but trained on the four general classes (Car, Cat, Dog, Flower). The first baseline assigns each detection of a general class to one of its subclasses sampled uniformly at random. The second baseline instead assigns each general detection to its the most frequent subclass. Finally, the third baseline assigns each general detection to a random subclass sampled according to their prior probabilities (as observed on the training set). Note that our second and third baselines require statistics of the subclasses, which are not available when considering only the box-level classes.

*Results.* Results are presented in Table 17. While all baselines yield poor results below < 0.05 mAP, our method delivers decent fine-grained detectors, with mAP ranging from 0.231 over the 61 subclasses of Cat, to 0.594 over the 102 subclasses of Flower. Interestingly, for several subclasses we have very good results suggesting that these classes are very distinctive, e.g. Ferrari (0.638 mAP), Land Rover (0.620 mAP), and Schnauzer (0.542 mAP). Several example detections are shown in Figure 31. These results demonstrate that the unified annotations of Open Images enable to train object detectors for fine-grained classes despite having only bounding box annotations for their parent class.

### 6.2 Zero-shot visual relationship detection by combining object bounding boxes and visual relationships

In the classical zero-shot Visual Relationship Detection (VRD) task, the zero-shot triplets consist only of new combinations of classes appearing in other annotated relationships (Lu et al., 2016; Liang et al., 2018), e.g. detecting ⟨ Cat, under, Table ⟩ when the training set contains ⟨ Cat, behind, Door ⟩ and ⟨ Person, under, Table ⟩. We propose to detect relationships also for new classes that are not present in any relationship annotations, by leveraging the bounding

| Bounding-box class `cls` | Number of subclasses | Number of samples | Examples of frequent and infrequent subclasses `C(cls)` (number of samples in parentheses) |
|---|---|---|---|
| `Car` | 62 | 4254 | `Ford` (354), `Chevrolet` (223), ..., `Frazer Nash` (2), `Riley Motor` (1) |
| `Cat` | 77 | 1340 | `Arabian Mau` (62), `American Wirehair` (56), ..., `Donskoy` (1), `Minskin` (1) |
| `Dog` | 39 | 4405 | `Terrier` (245), `Pinscher` (166), ..., `Malshi` (2), `Beaglier` (1) |
| `Flower` | 218 | 2541 | `Orchids` (119), `Buttercups` (114), ..., `Tidy tips` (1), `Aechmea 'Blue Tango'` (1) |

**Table 16** Statistics of the more specific classes for `Car`, `Flower`, `Cat`, and `Dog`. In all cases we report the total number of samples over the training and test sets.

| General class `cls` | Num. of fine-grained classes ($\geq$ 5 samples) | Uniform random sub-class (mAP) | Most common sub-class (mAP) | Prior-based random sub-class (mAP) | Image label transfer (ours) (mAP) |
|---|---|---|---|---|---|
| `Car` | 57 | 0.008 | 0.002 | 0.011 | 0.287 |
| `Cat` | 61 | 0.010 | 0.002 | 0.041 | 0.231 |
| `Dog` | 33 | 0.018 | 0.011 | 0.010 | 0.272 |
| `Flower` | 102 | 0.002 | 0.000 | 0.009 | 0.594 |

**Table 17** Results on fine-grained detection over subclasses of `Car`, `Flower`, `Cat`, and `Dog`.
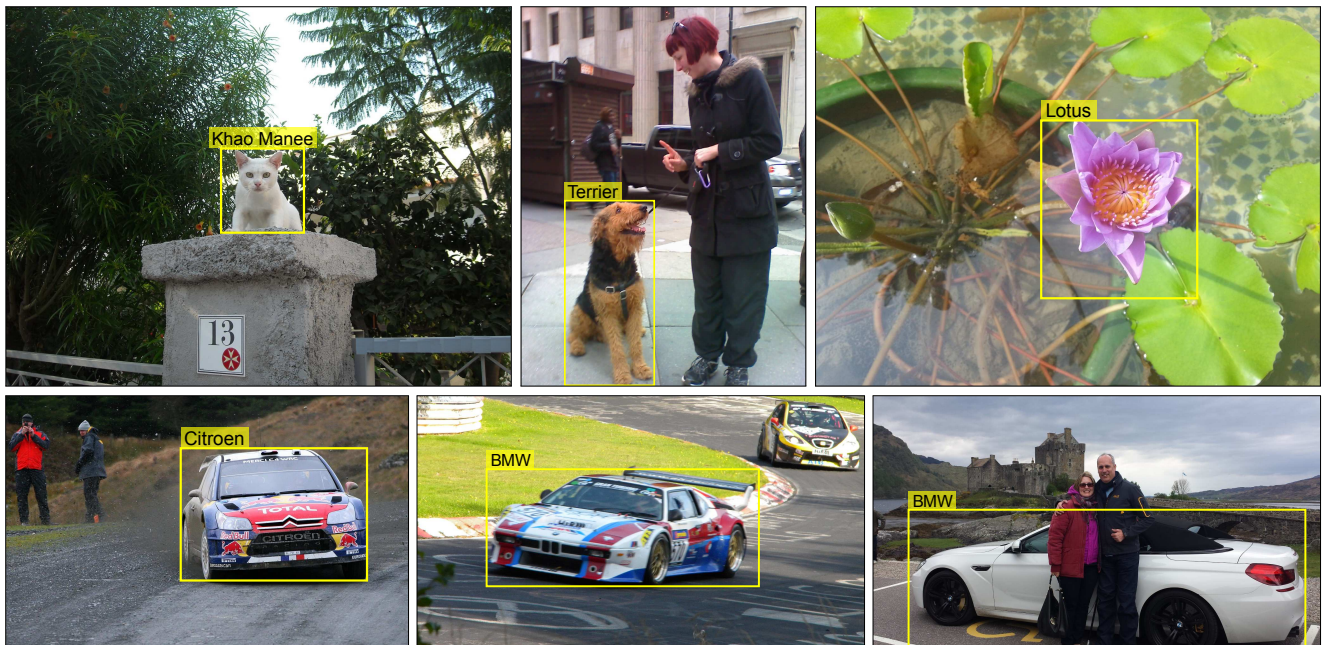


**Fig. 31** Example output of our fine-grained detectors. Note that our detector correctly identifies the individual car brands even in race cars whose appearance has been heavily modified (e.g. two bottom left examples).

box annotations that co-exists in the same images as relationship annotations. This showcases the benefits of the unified annotations in Open Images. In particular, we combine training data for related triplets that have the same subject and relationship with bounding boxes for a new object, e.g. detecting ⟨`Cat`, under, `Car`⟩ when we only have ⟨`Cat`, under, `Table`⟩ annotated as relationship and `Car` as bounding boxes. Analogously, we also introduce new subjects.

Specifically, we select new classes that are among the 600 with boxes in Open Images and do not have relationship annotations, and use them as our zero-shot classes. For evaluation purposes, we annotated the additional visual relationships for the new classes on the test split. We then fil-

tered the set of new zero-shot classes as those that are frequent enough (there are at least 10 instances of visual relationships in the test split that contain that class). Overall, we obtain 194 new triplets with a relationship preposition existing in the training split, where either subject or object are from the set of the new zero-shot classes. As a result, we have 5,983 new zero-shot triplet annotations on the test split, covering 48 new zero-shot classes, from which we have 284k annotated bounding boxes in the training split. Table 19 shows examples of the new object classes and the relationship triplets they are involved with.

We train the BAR-CNN model (Kolesnikov et al., 2018) using the annotated visual relationships on the training set

| Metric | R@50 (all triplets) | R@100 (all triplets) | R@50 (zero-shot) | R@100 (zero-shot) |
|---|---|---|---|---|
| Relationship detection | 40.61 | 40.93 | 7.68 | 7.70 |
| Phrase detection | 43.65 | 43.86 | 10.98 | 11.08 |

**Table 18 Zero-shot visual relationship detection results**. We evaluate a single model trained on the existing VRD annotations and box annotations. We report performance on all test set annotations including supervised and zero-shot triplets, as well as the performance on zero-shot-only triplets in terms of recall for relationship and predicate detection (Lu et al., 2016).
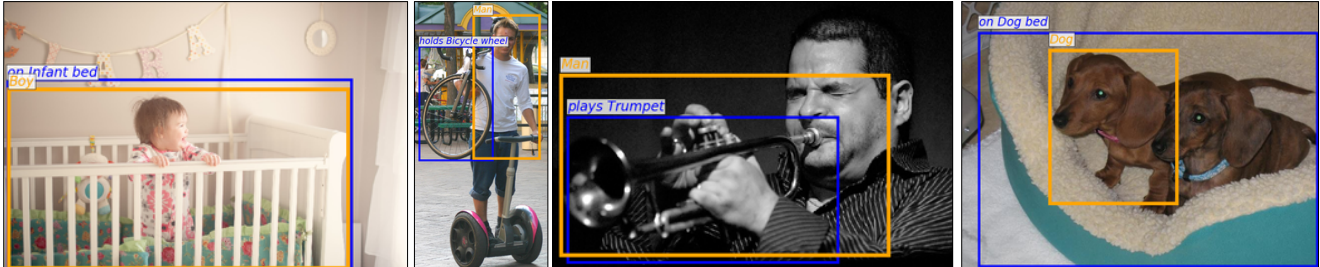


**Fig. 32 Example zero-shot detections** on the Open Images test set. Yellow boxes denote the subjects and blue boxes denote the objects and relationships. For these examples, object classes are zero-shot.

| Relationship | Zero-shot triplets |
|---|---|
| plays (26) | ⟨Girl, plays, **Cello**⟩, ⟨Man, plays, **Saxophone**⟩, . . . |
| holds (70) | ⟨Man, holds, **Bicycle wheel**⟩, ⟨Boy, holds, **Skateboard**⟩, . . . |
| wears (76) | ⟨Woman, wears, **Scarf**⟩, ⟨Man, wears, **Glasses**⟩, . . . , ⟨Girl, wears, **Necklace**⟩ |
| on (10) | ⟨Spoon, on, **Cutting board**⟩, ⟨Dog, on, **Dog bed**⟩, . . . |
| inside of (7) | ⟨Woman, inside of, **Golf cart**⟩, ⟨Girl, inside of, **Bus**⟩, . . . |
| at (4) | ⟨Boy, at, **Billiard Table**⟩, ⟨Girl, at, **Billiard Table**⟩, ⟨Woman, at, **Billiard Table**⟩, ⟨Man, at, **Billiard Table**⟩ |
| under (1) | ⟨Cat, under, **Coffee Table**⟩ |

**Table 19 Examples of zero-shot relationship triplets**, involving the new classes in bold. In parentheses, the total number of zero-shot triplets for each relationship.
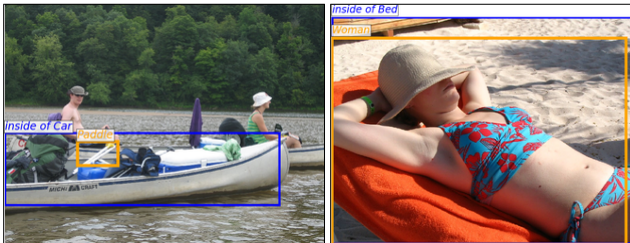


**Fig. 33 Interesting wrong predictions of the zero-shot model**. **Left image**: BAR-CNN detects the subject and relationship preposition correctly but the object label is incorrect (Boat is a zero-shot class). **Right image**: BAR-CNN only detects the subject correctly but both the relationship and the object are very close in meaning to what is shown in the image.

and the additional set of 48 classes for which only box annotations are available. During training, BAR-CNN accepts two types of samples: train samples with a single subject class label and train samples with object and relationship class labels. A combination of classification and box regression loss is optimized, as described in (Lin et al., 2017). In BAR-CNN, a sigmoid cross-entropy loss is used to handle multi-class samples with object and relationship class labels.

Since for the zero-shot object samples the relationship class labels are not available, we mask out the sigmoid loss components for the object samples without relationship class labels (those are samples derived from box annotations).

Table 18 presents our quantitative evaluation. The results show that the highly challenging setting of zero-shot VRD can be tackled to a reasonable degree but the gap to the supervised results is still very large, indicating potential for further exploration of the task with help of the unified annotations of Open Images. Figure 32 shows some examples of BAR-CNN zero-shot predictions and Figure 33 shows interesting prediction examples: incorrectly detected object class and incorrectly detected relationship (according to defined triplets); note that in both cases general semantics is preserved.

## 7 Conclusions

This paper presented Open Images V4, a collection of 9.2 million images annotated with unified ground-truth for image classification, object detection, and visual relationship detection. We explained how the data was collected and annotated, we presented comprehensive dataset statistics, we evaluated its quality, and we reported the performance of several modern models for image classification and object detection. We hope that the scale, quality, and variety of Open Images V4 will foster further research and innovation even beyond the areas of image classification, object detection, and visual relationship detection.

# 8 Credits

This is a full list of the contributors to the Open Images Dataset, which goes beyond the authors of this paper. Thanks to everyone!

*Project Lead and Coordination*: Vittorio Ferrari, Tom Duerig, and Victor Gomes.

*Image collection*: Ivan Krasin, David Cai.

*Image-level labels*: Neil Alldrin, Ivan Krasin, Shahab Kamali, Tom Duerig, Zheyun Feng, Anurag Batra, Alok Gunjan.

*Bounding boxes*: Hassan Rom, Alina Kuznetsova, Jasper Uijlings, Stefan Popov, Matteo Malloci, Sami Abu-El-Haija, Vittorio Ferrari.

*Visual relationships*: Alina Kuznetsova, Matteo Malloci, Vittorio Ferrari.

*Website and visualizer*: Jordi Pont-Tuset.

*Classes and hierarchy*: Chen Sun, Kevin Murphy, Tom Duerig, Vittorio Ferrari.

*Challenge*: Vittorio Ferrari, Alina Kuznetsova, Jordi Pont-Tuset, Matteo Malloci, Jasper Uijlings, Jake Walker, Rodrigo Benenson.

*Advisers*: Andreas Veit, Serge Belongie, Abhinav Gupta, Dhyanesh Narayanan, Gal Chechik.

# References

Alexe B, Deselaers T, Ferrari V (2010) What is an object? In: CVPR

Alexe B, Deselaers T, Ferrari V (2012) Measuring the objectness of image windows. IEEE Trans on PAMI

Chollet F (2017) Xception: Deep learning with depthwise separable convolutions. In: CVPR

Dai B, Zhang Y, Lin D (2017) Detecting visual relationships with deep relational networks. In: CVPR

Deng J, Dong W, Socher R, Li LJ, Li K, Fei-fei L (2009) ImageNet: A large-scale hierarchical image database. In: CVPR

Everingham M, Van Gool L, Williams CKI, Winn J, Zisserman A (2010) The PASCAL Visual Object Classes (VOC) Challenge. IJCV

Everingham M, Van Gool L, Williams CKI, Winn J, Zisserman A (2012) The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results. http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html

Everingham M, Eslami S, van Gool L, Williams C, Winn J, Zisserman A (2015) The PASCAL visual object classes challenge: A retrospective. IJCV

Fei-Fei L, Fergus R, Perona P (2006) One-shot learning of object categories. IEEE Trans on PAMI 28(4):594–611

Felzenszwalb P, Girshick R, McAllester D (2010a) Cascade Object Detection with Deformable Part Models. In: CVPR

Felzenszwalb P, Girshick R, McAllester D, Ramanan D (2010b) Object detection with discriminatively trained part based models. IEEE Trans on PAMI 32(9)

Gao C, Zou Y, Huang JB (2018) iCAN: Instance-centric attention network for human-object interaction detection. In: BMVC

Girshick R (2015) Fast R-CNN. In: ICCV

Girshick R, Donahue J, Darrell T, Malik J (2014) Rich feature hierarchies for accurate object detection and semantic segmentation. In: CVPR

Gkioxari G, Girshick R, Dollár P, He K (2018) Detecting and recognizing human-object interactions. CVPR

Griffin G, Holub A, Perona P (2007) The Caltech-256. Tech. rep., Caltech

Gupta A, Kembhavi A, Davis L (2009) Observing human-object interactions: Using spatial and functional compatibility for recognition. In: IEEE Trans. on PAMI

Gupta S, Malik J (2015) Visual semantic role labeling. arXiv preprint arXiv:150504474

He K, Zhang X, Ren S, Sun J (2016) Deep residual learning for image recognition. In: CVPR

Hinton GE, Vinyals O, Dean J (2014) Distilling the knowledge in a neural network. In: NeurIPS

Huang J, Rathod V, Sun C, Zhu M, Korattikara A, Fathi A, Fischer I, Wojna Z, Song Y, Guadarrama S, Murphy K (2017) Speed/accuracy trade-offs for modern convolutional object detectors. In: CVPR

Ioffe S, Szegedy C (2015) Batch normalization: Accelerating deep network training by reducing internal covariate shift. In: ICML

Kolesnikov A, Kuznetsova A, Lampert C, Ferrari V (2018) Detecting visual relationships using box attention. arXiv 1807.02136

Krishna R, Zhu Y, Groth O, Johnson J, Hata K, Kravitz J, Chen S, Kalantidis Y, Li LJ, Shamma DA, Bernstein M, Fei-Fei L (2017) Visual genome: Connecting language and vision using crowdsourced dense image annotations. IJCV 123(1):32–73

Krizhevsky A (2009) Learning multiple layers of features from tiny images. Tech. rep., University of Toronto

Krizhevsky A, Sutskever I, Hinton GE (2012) Imagenet classification with deep convolutional neural networks. In: NeurIPS

Li Y, Ouyang W, Wang X, Tang X (2017) ViP-CNN: Visual phrase guided convolutional neural network. In: CVPR

Liang K, Guo Y, Chang H, Chen X (2018) Visual relationship detection with deep structural ranking. In: AAAI

Liang X, Lee L, Xing EP (2017) Deep variation-structured reinforcement learning for visual relationship and attribute detection. In: CVPR

Lin T, Goyal P, Girshick R, He K, Dollar P (2017) Focal loss for dense object detection. In: ICCV

Lin TY, Maire M, Belongie S, Bourdev L, Girshick R, Hays J, Perona P, Ramanan D, Zitnick CL, Dollár P (2014) Microsoft COCO: Common objects in context. In: ECCV

Liu W, Anguelov D, Erhan D, Szegedy C, Reed S, Fu CY, Berg AC (2016) SSD: Single shot multibox detector. In: ECCV

Lu C, Krishna R, Bernstein M, Fei-Fei L (2016) Visual relationship detection with language priors. In: European Conference on Computer Vision

Mikolov T, Sutskever I, Chen K, Corrado GS, Dean J (2013) Distributed representations of words and phrases and their compositionality. In: NeurIPS

Papadopoulos DP, Uijlings JRR, Keller F, Ferrari V (2016) We don't need no bounding-boxes: Training object class detectors using only human verification. In: CVPR

Papadopoulos DP, Uijlings JR, Keller F, Ferrari V (2017) Extreme clicking for efficient object annotation. In: ICCV

Peyre J, Laptev I, Schmid C, Sivic J (2017) Weakly-supervised learning of visual relations. In: CVPR

Prest A, Schmid C, Ferrari V (2012) Weakly supervised learning of interactions between humans and objects. IEEE Trans on PAMI

Qian N (1999) On the momentum term in gradient descent learning algorithms. Neural Networks 12(1):145–151

Redmon J, Farhadi A (2017) YOLO9000: better, faster, stronger. In: CVPR

Redmon J, Divvala S, Girshick R, Farhadi A (2016) You only look once: Unified, real-time object detection. In: CVPR

Ren S, He K, Girshick R, Sun J (2015) Faster R-CNN: Towards real-time object detection with region proposal networks. In: NeurIPS

Russakovsky O, Deng J, Su H, Krause J, Satheesh S, Ma S, Huang Z, Karpathy A, Khosla A, Bernstein M, Berg A, Fei-Fei L (2015) ImageNet large scale visual recognition challenge. IJCV

Sandler M, Howard AG, Zhu M, Zhmoginov A, Chen L (2018) Mobilenetv2: Inverted residuals and linear bottleneck. In: CVPR

Su H, Deng J, Fei-Fei L (2012) Crowdsourcing annotations for visual object detection. In: AAAI Human Computation Workshop

Sun C, Shrivastava A, Singh S, Gupta A (2017) Revisiting unreasonable effectiveness of data in deep learning era. In: ICCV

Szegedy C, Liu W, Jia Y, Sermanet P, Reed S, Anguelov D, Erhan D, Vanhoucke V, Rabinovich A (2015) Going deeper with convolutions. In: CVPR

Szegedy C, Vanhoucke V, Ioffe S, Shlens J, Wojna Z (2016) Rethinking the inception architecture for computer vision. In: CVPR

Szegedy C, Ioffe S, Vanhoucke V, Alemi A (2017) Inception-v4, inception-resnet and the impact of residual connections on learning. In: AAAI

Uijlings J, Popov S, Ferrari V (2018) Revisiting knowledge transfer for training object class detectors. In: CVPR

Uijlings JRR, van de Sande KEA, Gevers T, Smeulders AWM (2013) Selective search for object recognition. IJCV

Veit A, Alldrin N, Chechik G, Krasin I, Gupta A, Belongie S (2017) Learning from noisy large-scale datasets with minimal supervision. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp 839–847, URL http://openaccess.thecvf.com/content_cvpr_2017/papers/Veit_Learning_From_Noisy_CVPR_2017_paper.pdf

Viola P, Jones M (2001a) Rapid object detection using a boosted cascade of simple features. In: CVPR

Viola P, Jones M (2001b) Robust real-time object detection. IJCV

Xu D, Zhu Y, Choy C, Fei-Fei L (2017) Scene graph generation by iterative message passing. In: Computer Vision and Pattern Recognition (CVPR)

Yao B, Fei-Fei L (2010) Modeling mutual context of object and human pose in human-object interaction activities. In: CVPR

Zellers R, Yatskar M, Thomson S, Choi Y (2018) Neural motifs: Scene graph parsing with global context. In: CVPR

Zhang H, Kyaw Z, Chang SF, Chua TS (2017a) Visual translation embedding network for visual relation detection. In: CVPR

Zhang H, Kyaw Z, Yu J, Chang SF (2017b) PPR-FCN: weakly supervised visual relation detection via parallel pairwise R-FCN. In: ICCV