

# LVIS: A Dataset for Large Vocabulary Instance Segmentation

Agrim Gupta Piotr Dollár Ross Girshick

Facebook AI Research (FAIR)

## Abstract

*Progress on object detection is enabled by datasets that focus the research community’s attention on open challenges. This process led us from simple images to complex scenes and from bounding boxes to segmentation masks. In this work, we introduce LVIS (pronounced ‘el-vis’): a new dataset for Large Vocabulary Instance Segmentation. We plan to collect ~2 million high-quality instance segmentation masks for over 1000 entry-level object categories in 164k images. Due to the Zipfian distribution of categories in natural images, LVIS naturally has a long tail of categories with few training samples. Given that state-of-the-art deep learning methods for object detection perform poorly in the low-sample regime, we believe that our dataset poses an important and exciting new scientific challenge. LVIS is available at <http://www.lvisdataset.org>.*

## 1. Introduction

A central goal of computer vision is to endow algorithms with the ability to intelligently describe images. Object detection is a canonical image description task; it is intuitively appealing, useful in applications, and straightforward to benchmark in existing settings. The accuracy of object detectors has improved dramatically and new capabilities, such as predicting segmentation masks and 3D representations, have been developed. There are now exciting opportunities to push these methods towards new goals.

Today, rigorous evaluation of general purpose object detectors is mostly performed in the few category regime (*e.g.* 80) or when there are a large number of training examples per category (*e.g.* 100 to 1000+). Thus, there is an opportunity to enable research in the natural setting where there are a large number of categories *and* per-category data is sometimes scarce. *The long tail of rare categories is inescapable; annotating more images simply uncovers previously unseen, rare categories (see Fig. 9 and [38, 33, 31, 35]).* Efficiently learning from few examples is a significant open problem in machine learning and computer vision, making this opportunity one of the most exciting from a scientific and practical perspective. But to open this area to empirical study, a suitable, high-quality dataset and benchmark is required.



**Figure 1. Example annotations.** We present **LVIS**, a new dataset for benchmarking Large Vocabulary Instance Segmentation in the 1000+ category regime with a challenging long tail of rare objects.

We aim to enable this new research direction by designing and collecting **LVIS** (pronounced ‘el-vis’)—a benchmark dataset for research on Large Vocabulary Instance Segmentation. We are collecting instance segmentation masks for more than 1000 entry-level object categories (see Fig. 1). When completed, we plan for our dataset to contain 164k images and ~2 million *high-quality* instance masks.<sup>1</sup> Our annotation pipeline starts from a set of images that were collected without prior knowledge of the categories that will be labeled in them. We engage annotators in an iterative object spotting process that uncovers the long tail of categories that naturally appears in the images and avoids using machine learning algorithms to automate data labeling.

We designed a crowdsourced annotation pipeline that enables the collection of our large-scale dataset while also yielding high-quality segmentation masks. Quality is important for future research because relatively coarse masks, such as those in the COCO dataset [23], limit the ability to differentiate algorithm-predicted mask quality beyond a certain, coarse point. When compared to expert annotators, our segmentation masks have higher overlap and boundary

<sup>1</sup>We plan to annotate the 164k images in COCO 2017 (we have permission to label test2017); ~2M is a projection after labeling 85k images.

consistency than both COCO and ADE20K [37].

To build our dataset, we adopt an *evaluation-first design principle*. This principle states that we should first determine exactly how to perform quantitative evaluation and only then design and build a dataset collection pipeline to gather the data entailed by the evaluation. We select our benchmark task to be COCO-style instance segmentation and we use the same COCO-style average precision (AP) metric that averages over categories and different mask intersection over union (IoU) thresholds [24]. Task and metric continuity with COCO reduces barriers to entry.

Buried within this seemingly innocuous task choice are immediate technical challenges: How do we fairly evaluate detectors when one object can reasonably be labeled with multiple categories (see Fig. 2)? How do we make the annotation workload feasible when labeling 164k images with segmented objects from over 1000 categories?

The essential design choice resolving these challenges is to build a *federated dataset*: a single dataset that is formed by the union of a large number of smaller constituent datasets, each of which looks exactly like a traditional object detection dataset for a single category. Each small dataset provides the essential guarantee of *exhaustive annotations* for a single category—*all instances of that category are annotated*. Multiple constituent datasets may overlap and thus a single object within an image can be labeled with multiple categories. Furthermore, since the exhaustive annotation guarantee only holds within each small dataset, we do not require the entire federated dataset to be exhaustively annotated with all categories, which dramatically reduces the annotation workload. Crucially, at test time the membership of each image with respect to the constituent datasets is not known by the algorithm and thus it must make predictions as if all categories will be evaluated. The evaluation oracle evaluates each category fairly on its constituent dataset.

In the remainder of this paper, we summarize how our dataset and benchmark relate to prior work, provide details on the evaluation protocol, describe how we collected data, and then discuss results of the analysis of this data.

**Dataset Timeline.** We report detailed analysis on the 5000 image `val` subset that we have annotated twice. We have now annotated an additional 77k images (split between `train`, `val`, and `test`), representing ~50% of the final dataset; we refer to this as **LVIS v0.5** (see §A for details). The first LVIS Challenge, based on v0.5, will be held at the COCO Workshop at ICCV 2019.

## 1.1. Related Datasets

Datasets shape the technical problems researchers study and consequently the path of scientific discovery [21]. We owe much of our current success in image recognition to pioneering datasets such as MNIST [20], BSDS [26],

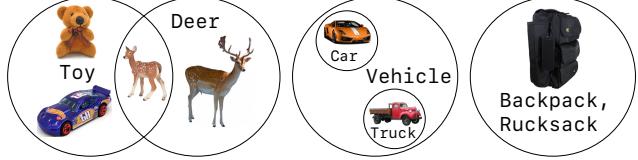


Figure 2. **Category relationships from left to right:** non-disjoint category pairs may be in *partially overlapping*, *parent-child*, or *equivalent (synonym)* relationships, implying that a single object may have multiple valid labels. The fair evaluation of an object detector must take the issue of multiple valid labels into account.

Caltech 101 [7], PASCAL VOC [6], ImageNet [30], and COCO [23]. These datasets enabled the development of algorithms that detect edges, perform large-scale image classification, and localize objects by bounding boxes and segmentation masks. They were also used in the discovery of important ideas, such as Convolutional Networks [19, 17], Residual Networks [13], and Batch Normalization [15].

LVIS is inspired by these and other related datasets, including those focused on street scenes (Cityscapes [4] and Mapillary [29]) and pedestrians (Caltech Pedestrians [5]). We review the most closely related datasets below.

**COCO** [23] is the most popular instance segmentation benchmark for common objects. It contains 80 categories that are pairwise distinct. There are a total of 118k training images, 5k validation images, and 41k test images. All 80 categories are exhaustively annotated in all images (ignoring annotation errors), leading to approximately 1.2 million instance segmentation masks. To establish continuity with COCO, we adopt the same instance segmentation task and AP metric, *and we are also annotating all images from the COCO 2017 dataset*. All 80 COCO categories can be mapped into our dataset. In addition to representing an order of magnitude more categories than COCO, our annotation pipeline leads to higher-quality segmentation masks that more closely follow object boundaries (see §4).

**ADE20K** [37] is an ambitious effort to annotate almost every pixel in 25k images with object instance, ‘stuff’, and part segmentations. The dataset includes approximately 3000 named objects, stuff regions, and parts. Notably, ADE20K was annotated by a *single expert annotator*, which increases consistency but also limits dataset size. Due to the relatively small number of annotated images, most of the categories do not have enough data to allow for both training and evaluation. Consequently, the instance segmentation benchmark associated with ADE20K evaluates algorithms on the 100 most frequent categories. In contrast, our goal is to enable benchmarking of *large vocabulary* instance segmentation methods.

**iNaturalist** [34] contains nearly 900k images annotated with bounding boxes for 5000 plant and animal species. Similar to our goals, iNaturalist emphasizes the importance

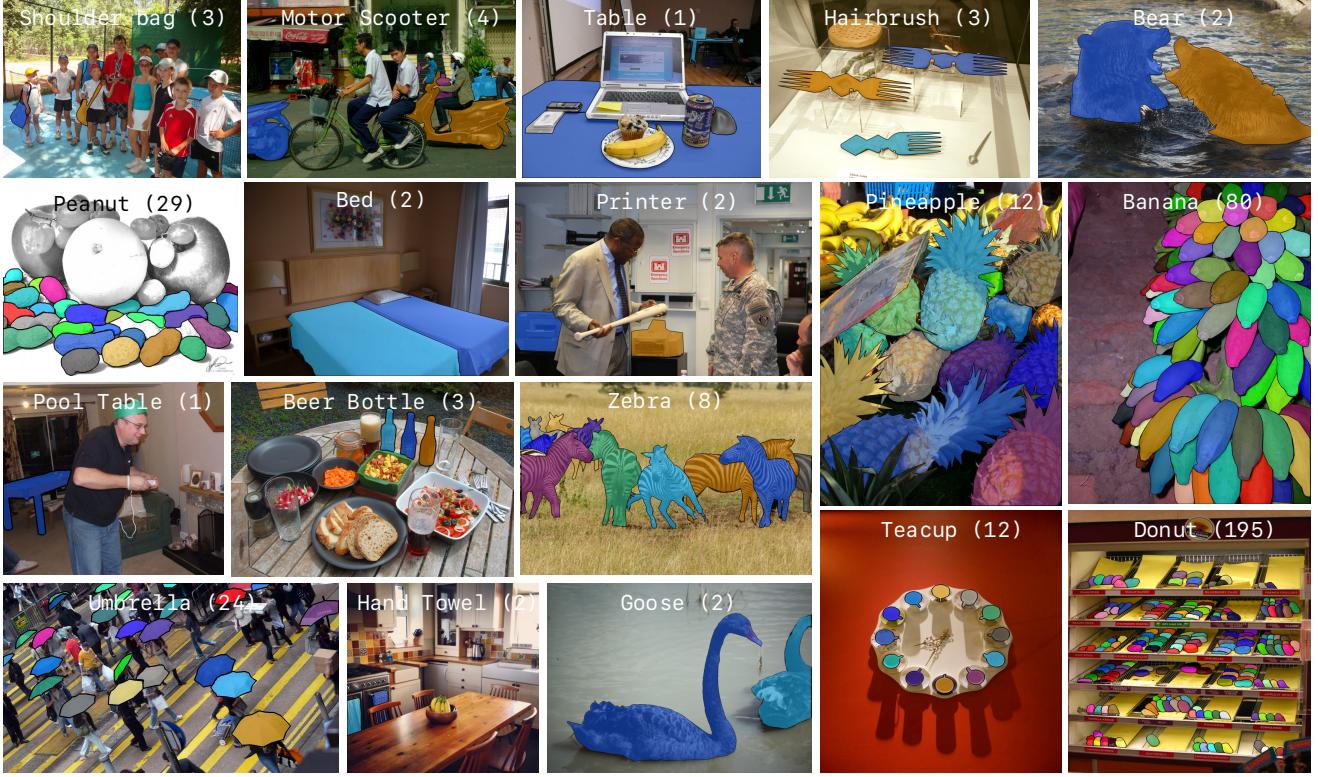


Figure 3. Example LVIS annotations (one category per image for clarity). See <http://www.lvisdataset.org/explore>.

of benchmarking classification and detection in the few example regime. Unlike our effort, iNaturalist does not include segmentation masks and is focussed on a different image and *fine-grained* category distribution; our category distribution emphasizes entry-level categories.

**Open Images v4** [18] is a large dataset of 1.9M images. The detection portion of the dataset includes 15M bounding boxes labeled with 600 object categories. The associated benchmark evaluates the 500 most frequent categories, all of which have over 100 training samples ( $>70\%$  of them have over 1000 training samples). Thus, unlike our benchmark, low-shot learning is not integral to Open Images. Also different from our dataset is the use of machine learning algorithms to select which images will be annotated by using classifiers for the target categories. Our data collection process, in contrast, involves no machine learning algorithms and instead discovers the objects that appear within a given set of images. Starting with release v4, Open Images has used a federated dataset design for object detection.

## 2. Dataset Design

We followed an *evaluation-first design principle*: prior to any data collection, we precisely defined what task would be performed and how it would be evaluated. This principle is important because there are technical challenges that arise when evaluating detectors on a large vocabulary dataset that

do not occur when there are few categories. These must be resolved first, because they have profound implications for the structure of the dataset, as we discuss next.

### 2.1. Task and Evaluation Overview

**Task and Metric.** Our dataset benchmark is the instance segmentation task: given a fixed, known set of categories, design an algorithm that when presented with a previously unseen image will output a segmentation mask for each instance of each category that appears in the image along with the category label and a confidence score. Given the output of an algorithm over a set of images, we compute *mask average precision* (AP) using the definition and implementation from the COCO dataset [24] (for more detail see §2.3).

**Evaluation Challenges.** Datasets like PASCAL VOC and COCO use manually selected categories that are *pairwise disjoint*: when annotating a *car*, there’s never any question if the object is instead a *potted plant* or a *sofa*. When increasing the number of categories, it is inevitable that other types of pairwise relationships will occur: (1) partially overlapping visual concepts; (2) parent-child relationships; and (3) perfect synonyms. See Fig. 2 for examples.

If these relations are not properly addressed, then the evaluation protocol will be unfair. For example, most *toys* are not *deer* and most *deer* are not *toys*, but a *toy deer* is both—if a detector outputs *deer* and the object is only labeled *toy*, the detection will be marked as wrong. Likewise,

if a car is only labeled *vehicle*, and the algorithm outputs *car*, it will be incorrectly judged to be wrong. Or, if an object is only labeled *backpack* and the algorithm outputs the synonym *rucksack*, it will be incorrectly penalized. Providing a fair benchmark is important for accurately reflecting algorithm performance.

These problems occur when the ground-truth annotations are missing one or more true labels for an object. If an algorithm happens to predict one of these correct, *but missing* labels, it will be unfairly penalized. Now, if all objects are exhaustively and correctly labeled with all categories, then the problem is trivially solved. But correctly and exhaustively labeling 164k images each with 1000 categories is undesirable: it forces a binary judgement deciding if each category applies to each object; there will be many cases of genuine ambiguity and inter-annotator disagreement. Moreover, the annotation workload will be very large. Given these drawbacks, we describe our solution next.

## 2.2. Federated Datasets

Our key observation is that the desired evaluation protocol does not require us to exhaustively annotate all images with all categories. What is required instead is that for each category  $c$  there must exist two disjoint subsets of the entire dataset  $\mathcal{D}$  for which the following guarantees hold:

**Positive set:** there exists a subset of images  $\mathcal{P}_c \subseteq \mathcal{D}$  such that all instances of  $c$  in  $\mathcal{P}_c$  are segmented. In other words,  $\mathcal{P}_c$  is exhaustively annotated for category  $c$ .

**Negative set:** there exists a subset of images  $\mathcal{N}_c \subseteq \mathcal{D}$  such that no instance of  $c$  appears in any of these images.

Given these two subsets for a category  $c$ ,  $\mathcal{P}_c \cup \mathcal{N}_c$  can be used to perform standard COCO-style AP evaluation for  $c$ . The evaluation oracle only judges the algorithm on a category  $c$  over the subset of images in which  $c$  has been exhaustively annotated; if a detector reports a detection of category  $c$  on an image  $i \notin \mathcal{P}_c \cup \mathcal{N}_c$ , the detection is *not* evaluated.

By collecting the per-category sets into a single dataset,  $\mathcal{D} = \cup_c (\mathcal{P}_c \cup \mathcal{N}_c)$ , we arrive at the concept of a *federated dataset*. A federated dataset is a dataset that is formed by the union of smaller constituent datasets, each of which looks exactly like a traditional object detection dataset for a single category. By not annotating all images with all categories, freedom is created to design an annotation process that avoids ambiguous cases and collects annotations only if there is sufficient inter-annotator agreement. At the same time, the workload can be dramatically reduced.

Finally, we note that positive set and negative set membership on the test split is not disclosed and therefore algorithms have no side information about what categories will be evaluated in each image. An algorithm thus must make its best prediction for *all* categories in each test image.

**Reduced Workload.** Federated dataset design allows us to make  $|\mathcal{P}_c \cup \mathcal{N}_c| \ll |\mathcal{D}|, \forall c$ . This choice dramatically re-

duces the workload and allows us to undersample the most frequent categories in order to avoid wasting annotation resources on them (*e.g.* *person* accounts for 30% of COCO). Of our estimated  $\sim 2$  million instances, likely no single category will account for more than  $\sim 3\%$  of the total instances.

## 2.3. Evaluation Details

The challenge evaluation server will only return the overall AP, not per-category AP’s. We do this because: (1) it avoids leaking which categories are present in the test set;<sup>2</sup> (2) given that tail categories are rare, there will be few examples for evaluation in some cases, which makes per-category AP unstable; (3) by averaging over a large number of categories, the overall category-averaged AP has lower variance, making it a robust metric for ranking algorithms.

**Non-Exhaustive Annotations.** We also collect an image-level boolean label,  $e_i^c$ , indicating if image  $i \in \mathcal{P}_c$  is exhaustively annotated for category  $c$ . In most cases (91%), this flag is true, indicating that the annotations are indeed exhaustive. In the remaining cases, there is at least one instance in the image that is not annotated. Missing annotations often occur in ‘crowds’ where there are a large number of instances and delineating them is difficult. During evaluation, we do not count false positives for category  $c$  on images  $i$  that have  $e_i^c$  set to false. We do measure recall on these images: the detector is expected to predict accurate segmentation masks for the labeled instances. Our strategy differs from other datasets that use a small maximum number of instances per image, per category (10-15) together with ‘crowd regions’ (COCO) or use a special ‘group of  $c$ ’ label to represent 5 or more instances (Open Images v4). Our annotation pipeline (§3) attempts to collect segmentations for *all* instances in an image, regardless of count, and then checks if the labeling is in fact exhaustive. See Fig. 3.

**Hierarchy.** During evaluation, we treat all categories the same; we do nothing special in the case of hierarchical relationships. To perform best, for each detected object  $o$ , the detector should output the most specific correct category as well as all more general categories, *e.g.*, a canoe should be labeled both *canoe* and *boat*. The detected object  $o$  in image  $i$  will be evaluated with respect to all labeled positive categories  $\{c \mid i \in \mathcal{P}_c\}$ , which may be any subset of categories between the most specific and the most general.

**Synonyms.** A federated dataset that separates synonyms into different categories is valid, but is unnecessarily fragmented (see Fig. 2, right). We avoid splitting synonyms into separate categories with WordNet [28]. Specifically, in LVIS each category  $c$  is a WordNet *synset*—a word sense specified by a set of synonyms and a definition.

<sup>2</sup>It’s possible that the categories present in the *val* and *test* sets may be a strict subset of those in the *train* set; we use the standard COCO 2017 *val* and *test* splits and cannot guarantee that all categories present in the *train* images are also present in *val* and *test*.

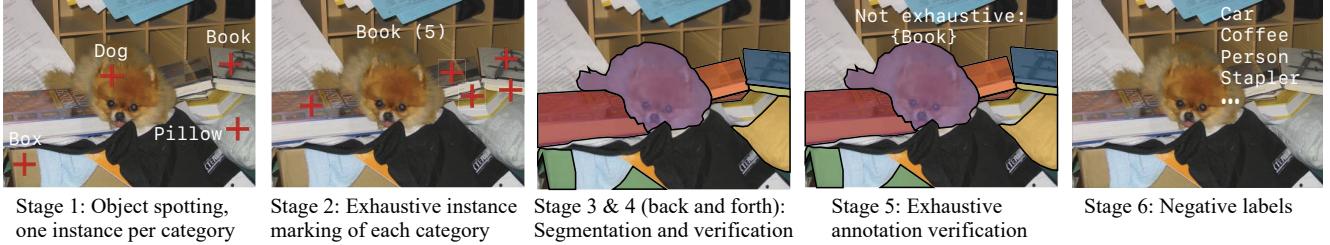


Figure 4. Our **annotation pipeline** comprises six stages. **Stage 1: Object Spotting** elicits annotators to mark a single instance of many different categories per image. This stage is iterative and causes annotators to discover a long tail of categories. **Stage 2: Exhaustive Instance Marking** extends the stage 1 annotations to cover all instances of each spotted category. Here we show additional instances of *book*. **Stages 3 and 4: Instance Segmentation and Verification** are repeated back and forth until  $\sim 99\%$  of all segmentations pass a quality check. **Stage 5: Exhaustive Annotations Verification** checks that all instances are in fact segmented and flags categories that are missing one or more instances. **Stage 6: Negative Labels** are assigned by verifying that a subset of categories do not appear in the image.

### 3. Dataset Construction

In this section we provide an overview of the annotation pipeline that we use to collect LVIS.

#### 3.1. Annotation Pipeline

Fig. 4 illustrates our annotation pipeline by showing the output of each stage, which we describe below. For now, assume that we have a fixed category vocabulary  $\mathcal{V}$ . We will describe how the vocabulary was collected in §3.2.

**Object Spotting, Stage 1.** The goals of the object spotting stage are to: (1) generate the positive set,  $\mathcal{P}_c$ , for each category  $c \in \mathcal{V}$  and (2) elicit vocabulary recall such that many different object categories are included in the dataset.

Object spotting is an iterative process in which each image is visited a variable number of times. On the first visit, an annotator is asked to mark one object with a point and to name it with a category  $c \in \mathcal{V}$  using an *autocomplete* text input. On each subsequent visit, all previously spotted objects are displayed and an annotator is asked to mark an object of a previously unmarked category or to skip the image if no more categories in  $\mathcal{V}$  can be spotted. When an image has been skipped 3 times, it will no longer be visited. The autocomplete is performed against the set of all synonyms, presented with their definitions; we internally map the selected word to its synset/category to resolve synonyms.

Obvious and salient objects are spotted early in this iterative process. As an image is visited more, less obvious objects are spotted, including incidental, non-salient ones. We run the spotting stage twice, and for each image we retain categories that were spotted in both runs. *Thus two people must independently agree on a name in order for it to be included in the dataset; this increases naming consistency.*

To summarize the output of stage 1: for each category in the vocabulary, we have a (possibly empty) set of images in which one object of that category is marked per image. This defines an initial positive set,  $\mathcal{P}_c$ , for each category  $c$ .

**Exhaustive Instance Marking, Stage 2.** The goals this stage are to: (1) verify stage 1 annotations and (2) take each image  $i \in \mathcal{P}_c$  and mark *all* instances of  $c$  in  $i$  with a point.

In this stage,  $(i, c)$  pairs from stage 1 are each sent to 5 annotators. They are asked to perform two steps. First, they are shown the definition of category  $c$  and asked to verify if it describes the spotted object. Second, if it matches, then the annotators are asked to mark all other instances of the same category. If it does not match, there is no second step. To prevent frequent categories from dominating the dataset and to reduce the overall workload, we subsample frequent categories such that no positive set exceeds more than 1% of the images in the dataset.

To ensure annotation quality, we embed a ‘gold set’ within the pool of work. These are cases for which we know the correct ground-truth. We use the gold set to automatically evaluate the work quality of each annotator so that we can direct work towards more reliable annotators. We use 5 annotators per  $(i, c)$  pair to help ensure instance-level recall.

To summarize, from stage 2 we have exhaustive instance spotting for each image  $i \in \mathcal{P}_c$  for each category  $c \in \mathcal{V}$ .

**Instance Segmentation, Stage 3.** The goals of the instance segmentation stage are to: (1) verify the category for each marked object from stage 2 and (2) upgrade each marked object from a point annotation to a full segmentation mask.

To do this, each pair  $(i, o)$  of image  $i$  and marked object instance  $o$  is presented to one annotator who is asked to verify that the category label for  $o$  is correct and if it is correct, to draw a *detailed* segmentation mask for it (e.g. see Fig. 3).

We use a training task to establish our quality standards. Annotator quality is assessed with a gold set and by tracking their average vertex count per polygon. We use these metrics to assign work to reliable annotators.

In sum, from stage 3 we have for each image and spotted instance pair one segmentation mask (if it is not rejected).

**Segment Verification, Stage 4.** The goal of the segment verification stage is to verify the quality of the segmentation masks from stage 3. We show each segmentation to

up to 5 annotators and ask them to rate its quality using a rubric. If two or more annotators reject the mask, then we requeue the instance for stage 3 segmentation. Thus we only accept a segmentation if 4 annotators agree it is high-quality. Unreliable workers from stage 3 are not invited to judge segmentations in stage 4; we also use rejections rates from this stage to monitor annotator reliability. We iterate between stages 3 & 4 a total of four times, each time only re-annotating rejected instances.

To summarize the output of stage 4 (after iterating back and forth with stage 3): we have a high-quality segmentation mask for  $>99\%$  of all marked objects.

**Full Recall Verification, Stage 5.** The full recall verification stage finalizes the positive sets. The goal is to find images  $i \in \mathcal{P}_c$  where  $c$  is not exhaustively annotated. We do this by asking annotators if there are any unsegmented instances of category  $c$  in  $i$ . We ask up to 5 annotators and require at least 4 to agree that annotation is exhaustive. As soon as two believe it is not, we mark the exhaustive annotation flag  $e_i^c$  as false. We use a gold set to maintain quality.

To summarize the output of stage 5: we have a boolean flag  $e_i^c$  for each image  $i \in \mathcal{P}_c$  indicating if category  $c$  is exhaustively annotated in image  $i$ . This finalizes the positive sets along with their instance segmentation annotations.

**Negative Sets, Stage 6.** The final stage of the pipeline is to collect a negative set  $\mathcal{N}_c$  for each category  $c$  in the vocabulary. We do this by randomly sampling images  $i \in \mathcal{D} \setminus \mathcal{P}_c$ , where  $\mathcal{D}$  is all images in the dataset. For each sampled image  $i$ , we ask up to 5 annotators if category  $c$  appears in image  $i$ . If any one annotator reports that it does, we reject the image. Otherwise  $i$  is added to  $\mathcal{N}_c$ . We sample until the negative set  $\mathcal{N}_c$  reaches a target size of 1% of the images in the dataset. We use a gold set to maintain quality.

To summarize, from stage 6 we have a negative image set  $\mathcal{N}_c$  for each category  $c \in \mathcal{V}$  such that the category does not appear in any of the images in  $\mathcal{N}_c$ .

### 3.2. Vocabulary Construction

We construct the vocabulary  $\mathcal{V}$  with an iterative process that starts from a large super-vocabulary and uses the object spotting process (stage 1) to winnow it down. We start from 8.8k synsets that were selected from WordNet by removing some obvious cases (*e.g.* proper nouns) and then finding the intersection with highly concrete common nouns [3]. This yields a high-recall set of concrete, and thus likely visual, entry-level synsets. We then apply object spotting to 10k COCO images with autocomplete against this super-vocabulary. This yields a reduced vocabulary with which we repeat the process once more. Finally, we perform minor manual editing. The resulting vocabulary contains 1723 synsets—the upper bound on the number of categories that can appear in LVIS.

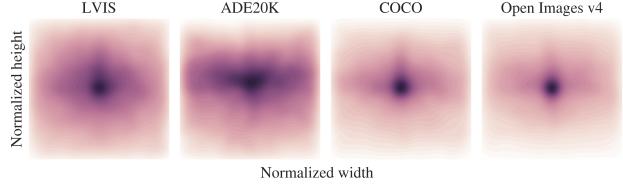


Figure 5. Distribution of object centers in normalized image coordinates for four datasets. ADE20K exhibits the greatest spatial diversity, with LVIS achieving greater complexity than COCO and the Open Images v4 training set.<sup>3</sup>

## 4. Dataset Analysis

For analysis, we have annotated 5000 images (the COCO `val2017` split) twice using the proposed pipeline. We begin by discussing general dataset statistics next before proceeding to an analysis of annotation consistency in §4.2 and an analysis of the evaluation protocol in §4.3.

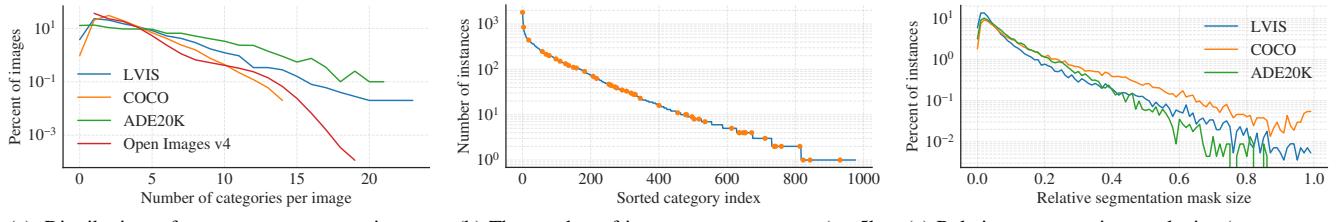
### 4.1. Dataset Statistics

**Category Statistics.** There are 977 categories present in the 5000 LVIS images. The category growth rate (see Fig. 9) indicates that the final dataset will have well over 1000 categories. On average, each image is annotated with 11.2 instances from 3.4 categories. The largest instances-per-image count is a remarkable 294. Fig. 6a shows the full categories-per-image distribution. LVIS’s distribution has more spread than COCO’s indicating that many images are labeled with more categories. The low-shot nature of our dataset can be seen in Fig. 6b, which plots the total number of instances for each category (in the 5000 images). The median value is 9, and while this number will be larger for the full image set, this statistic highlights the challenging long-tailed nature of our data.

**Spatial Statistics.** Our object spotting process (stage 1) encourages the inclusion of objects distributed throughout the image plane, not just the most salient foreground objects. The effect can be seen in Fig. 5 which shows object-center density plots. All datasets have some degree of center bias, with ADE20K and LVIS having the most diverse spatial distribution. COCO and Open Images v4 (training set<sup>3</sup>) have similar object-center distributions with a marginally lower degree of spatial diversity.

**Scale Statistics.** Objects in LVIS are also more likely to be small. Fig. 6c shows the relative size distribution of object masks: compared with COCO, LVIS objects tend to be smaller and there are fewer large objects (*e.g.*, objects that occupy most of an image are  $\sim 10\times$  less frequent). ADE20K has the fewest large objects overall and more medium ones.

<sup>3</sup>The CVPR 2019 version of this paper shows the distribution of the Open Images v4 *validation* set, which has more center bias. The peakiness is also exaggerated due to an intensity scaling artifact. For more details, see <https://storage.googleapis.com/openimages/web/factsfigures.html>.

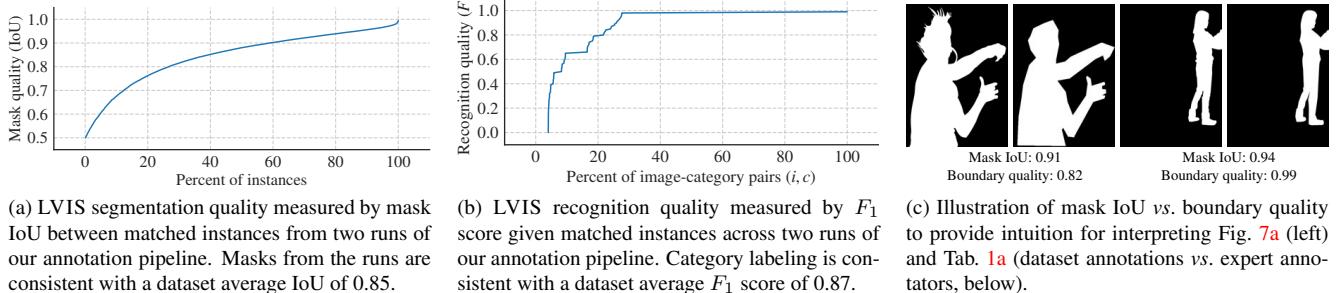


(a) Distribution of category count per image. LVIS has a heavier tail than COCO and Open Images training set. ADE20K is the most uniform.

(b) The number of instances per category (on 5k images) reveals the long tail with few examples. Orange dots: categories in common with COCO.

(c) Relative segmentation mask size (square root of mask-area-divided-by-image-area) compared between LVIS, COCO, and ADE20K.

Figure 6. **Dataset statistics.** Best viewed digitally.



(a) LVIS segmentation quality measured by mask IoU between matched instances from two runs of our annotation pipeline. Masks from the runs are consistent with a dataset average IoU of 0.85.

(b) LVIS recognition quality measured by  $F_1$  score given matched instances across two runs of our annotation pipeline. Category labeling is consistent with a dataset average  $F_1$  score of 0.87.

Figure 7. **Annotation consistency** using 5000 *doubly annotated* images from LVIS. Best viewed digitally.

dataset	comparison	mask IoU		boundary quality	
		mean	median	mean	median
COCO	dataset vs. experts	0.83 – 0.87	0.88 – 0.91	0.77 – 0.82	0.79 – 0.88
	expert 1 vs. expert 2	0.91 – 0.95	0.96 – 0.98	0.92 – 0.96	0.97 – 0.99
ADE20K	dataset vs. experts	0.84 – 0.88	0.90 – 0.93	0.83 – 0.87	0.84 – 0.92
	expert 1 vs. expert 2	0.90 – 0.94	0.95 – 0.97	0.90 – 0.95	0.99 – 1.00
LVIS	dataset vs. experts	<b>0.90 – 0.92</b>	<b>0.94 – 0.96</b>	<b>0.87 – 0.91</b>	<b>0.93 – 0.98</b>
	expert 1 vs. expert 2	0.93 – 0.96	0.96 – 0.98	0.91 – 0.96	0.97 – 1.00

(a) For each metric (mask IoU, boundary quality) and each statistic (mean, median), we show a bootstrapped 95% confidence interval. LVIS has the highest quality across all measures.

Table 1. Annotation quality and complexity relative to experts.

## 4.2. Annotation Consistency

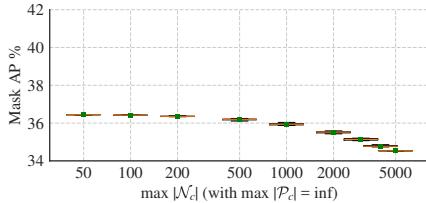
**Annotation Pipeline Repeatability.** A repeatable annotation pipeline implies that the process generating the ground-truth data is not overly random and therefore may be learned. To understand repeatability, we annotated the 5000 images twice: after completing object spotting (stage 1), we have initial positive sets  $\mathcal{P}_c$  for each category  $c$ ; we then execute stages 2 through 5 (exhaustive instance marking through full recall verification) twice in order to yield doubly annotated positive sets. To compare them, we compute a matching between them for each image and category pair. We find a matching that maximizes the total mask intersection over union (IoU) summed over the matched pairs and then discard any matches with  $\text{IoU} < 0.5$ . Given these matches we compute the dataset average mask IoU (0.85) and the dataset average  $F_1$  score (0.87). Intuitively, these quantities describe ‘segmentation quality’ and ‘recognition quality’ [16]. The cumulative distributions of these metrics (Fig. 7a and 7b) show that even though matches are estab-

dataset	annotation source	boundary complexity	
		mean	median
COCO	dataset	5.59 – 6.04	5.13 – 5.51
	experts	6.94 – 7.84	5.86 – 6.80
ADE20K	dataset	6.00 – 6.84	4.79 – 5.31
	experts	6.34 – 7.43	4.83 – 5.53
LVIS	dataset	<b>6.35 – 7.07</b>	<b>5.44 – 6.00</b>
	experts	7.13 – 8.48	5.91 – 6.82

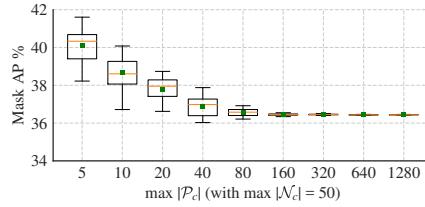
(b) Comparison of annotation complexity. Boundary complexity is perimeter divided by square root area [1].

lished based on a low IoU threshold (0.5), matched masks tend to have much higher IoU. The results show that roughly 50% of matched instances have IoU greater than 90% and roughly 75% of the image-category pairs have a perfect  $F_1$  score. Taken together, these metrics are a strong indication that our pipeline has a large degree of repeatability.

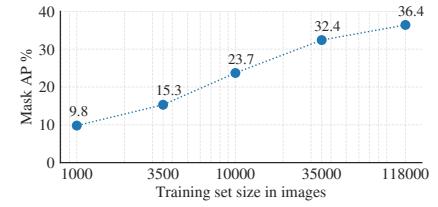
**Comparison with Expert Annotators.** To measure segmentation quality, we randomly selected 100 instances with mask area greater than  $32^2$  pixels from LVIS, COCO, and ADE20K. We presented these instances (indicated by bounding box and category) to two independent expert annotators and asked them to segment each object using professional image editing tools. We compare dataset annotations to expert annotations using mask IoU and boundary quality (boundary  $F$  [26]) in Tab. 1a. The results (bootstrapped 95% confidence intervals) show that our masks are high-quality, surpassing COCO and ADE20K on both measures (see Fig. 7c for intuition). At the same time, the objects in LVIS have more complex boundaries [1] (Tab. 1b).



(a) Given fixed detections, we show how AP varies with  $|\mathcal{N}_c|$ , the number of negative images per category used in evaluation.



(b) With the same detections from Fig. 8a and  $|\mathcal{N}_c| = 50$ , we show how AP varies as we vary  $|\mathcal{P}_c|$ , the positive set size.



(c) Low-shot detection is an open problem: training Mask R-CNN on 1k images decreases COCO val2017 mask AP from 36% to 10%.

Figure 8. Analysis of AP as a function of different data sizes. Best viewed digitally.

Mask R-CNN	test anno.	box AP	mask AP
ResNet-50-FPN model id: 35859007	COCO	38.2	34.1
	LVIS	38.8	34.4
ResNet-101-FPN model id: 35861858	COCO	40.6	36.0
	LVIS	40.9	36.0
ResNeXt-101-64x4d-FPN model id: 37129812	COCO	47.8	41.2
	LVIS	48.6	41.7

Table 2. COCO-trained Mask R-CNN evaluated on LVIS annotations. Both annotations yield similar AP values.

### 4.3. Evaluation Protocol

**COCO Detectors on LVIS.** To validate our annotations and federated dataset design we downloaded three Mask R-CNN [12] models from the Detectron Model Zoo [8] and evaluated them on LVIS annotations for the categories in COCO. Tab. 2 shows that both box AP and mask AP are close between our annotations and the original ones from COCO for all models, which span a wide AP range. This result validates our annotations and evaluation protocol: even though LVIS uses a federated dataset design with sparse annotations, the quantitative outcome closely reproduces the ‘gold standard’ results from dense COCO annotations.

**Federated Dataset Simulations.** For insight into how AP changes with positive and negative sets sizes  $|\mathcal{P}_c|$  and  $|\mathcal{N}_c|$ , we randomly sample smaller evaluation sets (20 times) from COCO val2017 and recompute AP. In Fig. 8a we use all positive instances for evaluation, but vary  $|\mathcal{N}_c|$  between 50 and 5k. AP decreases somewhat ( $\sim 2\%$  absolute) as we increase the number of negative images as the ratio of negative to positive examples grows with fixed  $|\mathcal{P}_c|$  and increasing  $|\mathcal{N}_c|$ . Next, in Fig. 8b we set  $|\mathcal{N}_c| = 50$  and vary  $|\mathcal{P}_c|$ . We observe that even with a small positive set size of 80, AP is similar to the baseline with low variance. With smaller positive sets (down to 5) variance increases, but the AP gap from 1st to 3rd quartile remains below 2% absolute. A curious upward bias in AP appears, which we investigate in §C.2. These simulations together with COCO detectors tested on LVIS (Tab. 2) indicate that including smaller evaluation sets for each category is viable for evaluation.

**Low-Shot Detection.** To validate the claim that low-shot detection is a challenging open problem, we trained Mask R-CNN on random subsets of COCO train2017 rang-

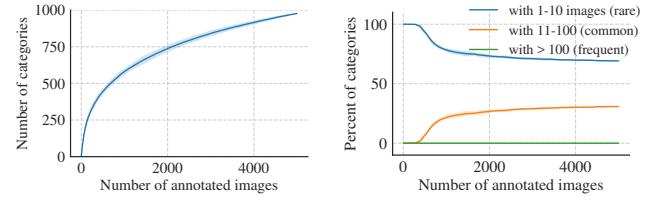


Figure 9. (Left) As more images are annotated, new categories are discovered. (Right) Consequently, the percentage of low-shot categories (blue curve) remains large, decreasing slowly.

ing from 1k to 118k images. For each subset, we optimized the learning rate schedule and weight decay by grid search. Results on val2017 are shown in Fig. 8c. At 1k images, mask AP drops from 36.4% (full dataset) to 9.8% (1k subset). In the 1k subset, 89% of the categories have more than 20 training instances, while the low-shot literature typically considers  $\ll 20$  examples per category [10].

**Low-Shot Category Statistics.** Fig. 9 (left) shows category growth as a function of image count (up to 977 categories in 5k images). Extrapolating the trajectory, our final dataset will include over 1k categories (upper bounded by the vocabulary size, 1723). Since the number of categories increases during data collection, the low-shot nature of LVIS is somewhat independent of the dataset scale, see Fig. 9 (right) where we bin categories based on how many images they appear in: *rare* (1-10 images), *common* (11-100), and *frequent* ( $> 100$ ). These bins, as measured w.r.t. the training set, will be used to present disaggregated AP metrics.

## 5. Conclusion

We introduced LVIS, a new dataset designed to enable, for the first time, the rigorous study of instance segmentation algorithms that can recognize a large vocabulary of object categories ( $> 1000$ ) and must do so using methods that can cope with the open problem of low-shot learning. While LVIS emphasizes learning from few examples, the dataset is not small: it will span 164k images and label  $\sim 2$  million object instances. Each object instance is segmented with a high-quality mask that surpasses the annotation quality of related datasets. We plan to establish LVIS as a benchmark challenge that we hope will lead to exciting new object detection, segmentation, and low-shot learning algorithms.

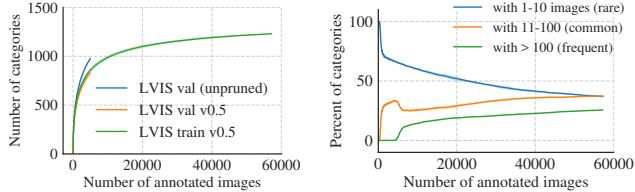


Figure 10. Category growth (left) and frequency statistics (right) for LVIS v0.5. Best viewed digitally. Compare with Fig. 9.

## A. LVIS Release v0.5

LVIS release v0.5 marks the halfway point in data collection. For this release, we have annotated an additional 77k images (57k train, 20k test) beyond the 5k val images that we analyzed in the previous sections, for a total of 82k annotated images. Release v0.5 is publicly available at <https://www.lvisdataset.org> and will be used in the first LVIS Challenge to be held in conjunction with the COCO Workshop at ICCV 2019.

**Collection Details.** We collected the v0.5 data in two 38.5k image batches using the process described in the main text. Each batch contained a proportional mix of train and test images. After collection was completed for the first batch, we manually checked all 1415 categories that were represented in the data collection and cast an include vs. exclude vote for each category based on its visual consistency. This process led to the removal of  $\sim 18\%$  of categories and  $\sim 10\%$  of labeled instances. After collecting the second batch, we repeated this process for 83 categories that were newly introduced. After we finish the full data collection for v1 (estimated early 2020), we will conduct another similar quality control pass on a subset of the categories.

LVIS val v0.5 is the same as the set used for analysis in the main text, except that we removed any categories that: (1) were determined to be visually inconsistent in the quality control pass or (2) had zero instances in the training set. In this section, we refer to the annotations used for analysis in the main text as ‘LVIS val (unpruned)’.

**Dataset Statistics.** After our quality control pass, the final category count for release v0.5 is 1230. The number of categories in the val set decreased from 977 to 830, due to quality control, and it now has 51k segmented object instances. The train v0.5 set has 694k segmented instances.

We next repeat some of the key analysis plots, this time showing the final val and train v0.5 sets compared to the original (unpruned) val set from the main text. The train and test sets are collected using an identical process (the train and test images were originally sampled from the same image distribution and are mixed together in each annotation batch) and therefore the training data is statistically identical to that of the test data.

Fig. 10 (left) illustrates the category growth rate on train and val before and after pruning. We expect

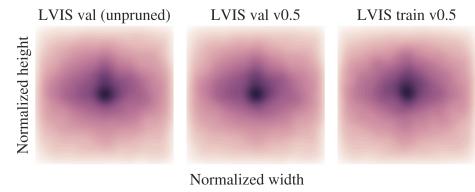


Figure 11. Distribution of object centers in normalized image coordinates for LVIS val (unpruned), LVIS val v0.5 (i.e. after quality control), and LVIS train v0.5. The distributions are nearly identical.

only modest growth while collecting the second half of the dataset, perhaps expanding by  $\sim 100$  additional categories. Next, we extend Fig. 9 (right) from 5k images to 57k images using the train v0.5 data, as shown in Fig. 10 (right). Due to the slowing category growth, the percent of rare categories (those appearing in 1-10 training images) is decreasing, but remains a sizeable portion of the dataset. Roughly 75% of categories appear in 100 training images or less, highlighting the challenging low-shot nature of the dataset.

Finally, we look at the spatial distribution of object centers in Fig. 11. This visualization verifies that quality control did not lead to a meaningful bias in this statistic. The train and val sets exhibit visually similar distributions.

Based on this analysis and our qualitative judgement when performing per-category quality control, we conclude that our data collection process scales well beyond the initial 5k set analyzed in the main text.

## B. LVIS v0.5 Baselines

To help researchers calibrate their results for the upcoming LVIS Challenge at ICCV 2019,<sup>4</sup> we introduce simple baselines. In §B.1, we test the performance of Mask R-CNN [12] out-of-the-box, and show the importance of adjusting two inference-time hyper-parameters. Next in §B.2 we provide an improved (yet standard) baseline that resamples the training data in order to increase the frequency of rare categories. Finally in §B.3 we train larger models.

### B.1. Mask R-CNN Out-of-the-Box

We first apply Mask R-CNN out-of-the-box on LVIS. Unless specified we use Mask R-CNN with a ResNet-50 backbone (pre-trained on ImageNet) with FPN [22]. Training is performed using Detectron2, which is implemented in PyTorch and will be open-sourced later this year. Our training formula is unmodified from COCO training.<sup>5</sup>

<sup>4</sup><https://www.lvisdataset.org/challenge>

<sup>5</sup>We use SGD with 0.9 momentum and 16 images per mini-batch; the training schedule is 60k/20k/10k updates at learning rates of 0.02/0.002/0.0002 respectively (this 90k update schedule is equivalent to  $\sim 25$  epochs over train v0.5); we use a linear learning rate warmup [9] over 1000 updates starting from a learning rate of 0.001; weight decay 0.0001 is applied; horizontal flipping is the only train-time data augmentation unless otherwise stated; training and inference images are resized to a shorter image edge of 800 pixels; no test-time augmentation is used.

score thr	det/img	AP	AP <sub>r</sub>	AP <sub>c</sub>	AP <sub>f</sub>	AP <sup>bb</sup>
0.050	100	14.8	0.8	10.9	25.3	14.8
0.050	300	15.7	0.8	12.1	26.1	15.6
0.001	300	20.8	3.3	20.7	27.9	20.3
0.000	300	20.9	3.4	20.9	27.9	20.4
0.050	100	14.8 $\pm$ 0.19	0.6 $\pm$ 0.21	11.0 $\pm$ 0.36	25.2 $\pm$ 0.10	14.8 $\pm$ 0.17
0.000	300	21.0 $\pm$ 0.17	3.2 $\pm$ 0.35	21.3 $\pm$ 0.45	27.7 $\pm$ 0.12	20.5 $\pm$ 0.21

(a) **Mask R-CNN baselines** (ResNet-50-FPN backbone). Top rows: adjusting two *inference-time* hyper-parameters, the minimum score threshold and the number of detections per image, leads to a gain of 6.1 AP over the baseline using standard COCO hyper-parameters (row 1). The last two rows show the mean and standard deviation from five training runs.

$t$	AP	AP <sub>r</sub>	AP <sub>c</sub>	AP <sub>f</sub>
0	21.0 $\pm$ 0.17	3.2 $\pm$ 0.35	21.3 $\pm$ 0.45	27.7 $\pm$ 0.12
0.0001	21.2 $\pm$ 0.14	4.5 $\pm$ 0.47	21.5 $\pm$ 0.37	27.6 $\pm$ 0.14
0.0010	23.2 $\pm$ 0.21	13.4 $\pm$ 0.80	23.2 $\pm$ 0.32	27.1 $\pm$ 0.07
0.0100	21.8 $\pm$ 0.25	9.8 $\pm$ 1.27	22.7 $\pm$ 0.48	25.6 $\pm$ 0.13
0.1000	21.3 $\pm$ 0.24	9.6 $\pm$ 0.83	21.7 $\pm$ 0.32	25.5 $\pm$ 0.10
CAS	18.7 $\pm$ 0.46	8.5 $\pm$ 1.56	19.0 $\pm$ 0.45	22.3 $\pm$ 0.19

(b) **Mask R-CNN with repeat factor sampling** (with best settings from Table 3a). The frequency threshold  $t$  controls the degree of resampling of rare categories ( $t=0$  gives no resampling). Setting  $t>0$  substantially improves AP<sub>r</sub> and  $t=0.001$  gives best overall results. The last row presents class aware sampling (CAS), an alternate oversampling method [32].

enhancement	AP	AP <sub>r</sub>	AP <sub>c</sub>	AP <sub>f</sub>
Table 3b best	23.2 $\pm$ 0.21	13.4 $\pm$ 0.80	23.2 $\pm$ 0.32	27.1 $\pm$ 0.07
+ scale jitter	24.4 $\pm$ 0.06	14.5 $\pm$ 0.67	24.3 $\pm$ 0.37	28.4 $\pm$ 0.12
+ ResNet-101	26.0 $\pm$ 0.18	15.8 $\pm$ 0.95	26.1 $\pm$ 0.21	29.8 $\pm$ 0.22
+ ResNeXt-101-32 $\times$ 8d	27.1 $\pm$ 0.43	15.6 $\pm$ 1.14	27.5 $\pm$ 0.77	31.4 $\pm$ 0.12

(c) **Mask R-CNN enhancements**. We apply scale jitter data augmentation and upgrade the backbone to larger models [36]. This improves all AP metrics although AP<sub>r</sub> does not improve with the largest backbone.

Table 3. **LVIS release v0.5 baselines**. Metrics: AP is mask AP; subscripts ‘r’, ‘c’, and ‘f’ refer to rare, common, and frequent category subsets (defined in §4.3). Where applicable we repeat each experiment 5 times and report mean and standard deviation.

The results are low and in particular AP<sub>r</sub> (mask AP for rare categories) is 0.8%—*near zero*. In Table 3a we demonstrate that adjusting two inference-time hyper-parameters on LVIS improves results. First, we increase the number of detections per images as LVIS allows up to 300 (vs. 100 for COCO). Second, due to class imbalance the max confidence scores reported for rare and common classes is typically low (compared to COCO), hence reducing the minimum score threshold from the default of 0.05 to 0.0 (*i.e.*, no threshold) substantially improves AP<sub>c</sub>. The combination of these changes increases AP<sub>r</sub> a modest amount, up to an average of 3.2% Table 3a (bottom row).

We additionally observe that on LVIS, *mask AP is typically slightly higher than box AP* (denoted by AP<sup>bb</sup>). This trend is the opposite for COCO, where AP is typically 3 to 4% (absolute) lower than AP<sup>bb</sup>. We hypothesize that the AP/AP<sup>bb</sup> trend on LVIS is due to high quality segmentation masks. We have found supporting evidence by programmatically degrading LVIS mask quality and observing a drop in AP with almost no change in AP<sup>bb</sup> (results not shown).

## B.2. Mask R-CNN with Data Resampling

Resampling training data is a common strategy for training models on class imbalanced datasets [32, 11, 25, 27]. We apply a method that was used to train large-scale hashtag prediction models in [25] (inspired by [27]). The method, which we refer to as *repeat factor sampling*, increases the rate at which tail categories are observed by oversampling the images that contain them.

The method is implemented as follows. For each category  $c$ , let  $f_c$  be the fraction of training images that contain at least one instance of  $c$ . Define the category-level repeat factor as  $r_c = \max(1, \sqrt{t/f_c})$ , where  $t$  is a hyper-parameter. Since each image may contain multiple categories, we define an image-level repeat factor. Specifically, for each image  $i$ , we set  $r_i = \max_{c \in i} r_c$ , where  $\{c \in i\}$  are the categories labeled in  $i$ . In each epoch, the SGD data sampler creates a random permutation of images in which each image is repeated according to its repeat factor  $r_i$ .

The one hyper-parameter of this method,  $t$ , is a threshold that intuitively controls the point at which oversampling kicks in. For categories with  $f_c \leq t$ , there is no oversampling. For categories with  $f_c > t$ , the degree of oversampling follows a square-root inverse frequency heuristic: if we decrease the frequency of a category by a factor  $\gamma < 1$ , then its repeat factor will be multiplied by  $\sqrt{1/\gamma}$ . This heuristic has worked well in other settings, *e.g.* [27].

The results of repeat factor sampling for varying  $t$  are shown in Table 3b. Comparing with the baseline (equivalent to  $t = 0$ ), there is a large improvement in AP<sub>r</sub> from 3.2% to 13.4% at  $t = 0.001$ . This threshold oversamples categories appearing in less than 0.1% of images (829 of the 1230 categories). There is a slight penalty in lower AP<sub>f</sub> ( $-0.6\%$ ), but overall AP improves (+2.2%).

We also present results using *class aware sampling* (CAS), a popular method on imbalanced classification datasets (*e.g.*, [32]). In CAS, the data sampler first selects a category and then an image containing that category. Consistent with repeat factor sampling and SGD best-practices [2], we iterate over random permutations of categories and within each category random permutations of their images. CAS improves AP<sub>r</sub> over the baseline as expected (from 3.2% to 8.5%), however both AP<sub>c</sub> and AP<sub>f</sub> decrease leading to a worse overall result.

## B.3. Mask R-CNN Standard Enhancements

Finally we consider two standard enhancements in addition to using repeat factor sampling with  $t = 0.001$ : we apply scale jitter at training time (sampling image scale from  $\{640, 672, 704, 736, 768, 800\}$ ) and upgrade to larger models. Both enhancement yield improvements as reported in Table 3c with a final validation AP of 27.1%.

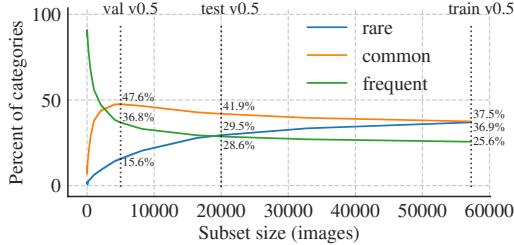


Figure 12. The distribution of rare, common, and frequent categories (defined w.r.t. `train v0.5`) within random image subsets of a given size changes as a function of that size. The shaded region (imperceptible without zoom) illustrates one standard deviation around the mean over 10 draws of subsets for each size.

## C. LVIS `val` to `test` Results Transfer

In this section we analyze how AP on LVIS v0.5 `val` transfers to the `test` set. First, in §C.1 we show how the category distribution varies between the smaller `val` and larger `test` sets. Next, in §C.2 we demonstrate the surprising result that even with a fixed category set *smaller evaluation sets can have a bias towards higher AP*. The impact is larger for rare categories, hence while it has a minimal effect on COCO, on LVIS it results in AP measured on the larger `test` set to be lower than on the smaller `val` set.

### C.1. Category Frequency Distributions

An *evaluation* set with a large proportion of categories that appear infrequently in the *training* set (*i.e.* categories with few training examples) will tend to be more difficult as learning from few examples is challenging. An important question, then, is what is the frequency distribution (w.r.t. the training set) of categories in a given evaluation set?

To investigate this question, we look at category frequency distributions in random image subsets of various sizes. For visualization, we quantize category frequency as described in §4.3 into ‘rare’, ‘common’, and ‘frequent’ groups based on how many images each category appears in the training set. In Fig. 12, for each subset size, we plot the mean rare/common/frequent category distribution over random subsets of that size.

From this analysis, we can predict that the `val` set (5k images) should contain  $\sim 15\%$  rare categories, while the `test` set (20k images) should have  $\sim 29\%$  rare categories. Their actual values are 15.1% and 28.2%, validating our prediction. In general, larger evaluation sets contain a higher proportion of rare categories. Using the  $AP_{r/c/f}$  achieved by R-101-FPN Mask R-CNN as an example, this distribution shift could result in a decrease in overall AP of  $\sim 2\%$  when moving from `val` to `test`. While optimizing for the LVIS Challenge, one may want to take this distribution shift into account; on the `test` set rare categories will play a more important role than on the `val` set.

### C.2. AP as a Function of Evaluation Set Size

Suppose we have a small evaluation set (*e.g.*, `val`, 5k images) and a large evaluation set (*e.g.*, `test`, 20k images) that are both random samples from the same population. One might expect that while the categories present in the sets are different, the *per-category AP* for a given category computed on evaluation sets of different sizes should be unbiased estimates of the true AP and only the variance of the estimate should change. Surprisingly, in general this intuition is not true and the estimate can be biased for smaller evaluation set sizes. We first observed this bias in Fig. 8b, in which we see that AP increases on average as the number of positive images per category ( $|\mathcal{P}_c|$ ) decreases. We now analyze this bias further in both simulated<sup>6</sup> and real data.

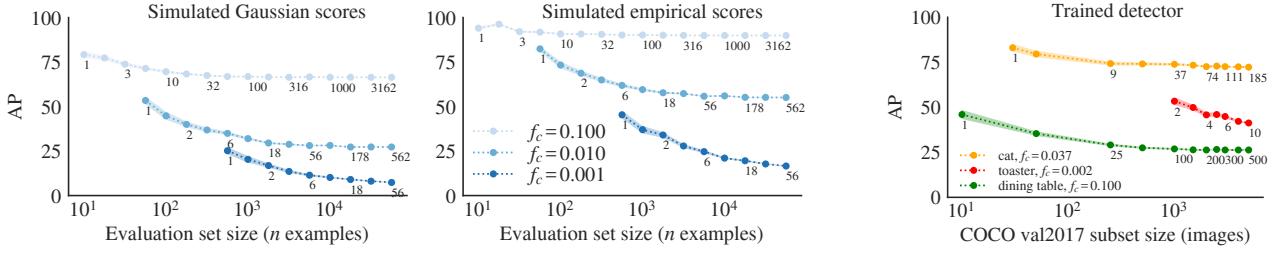
Consider a simple *binary classification* setting with  $n$  test examples in total. Each example is positive for category  $c$  with probability  $f_c$ . We will compute AP as a function of the evaluation set size  $n$  and for various values of  $f_c$ . We consider two simulated classifiers, each of which is defined by class-conditional probability distributions over scores:  $p_1(s) = p(s|y=1)$  and  $p_0(s) = p(s|y=0)$ , where  $s$  is the classifier score and  $y$  indicates the true label.

In the first simulation the classifier scores are drawn from  $p_1(s) = \mathcal{N}(1, 1)$  and  $p_0(s) = \mathcal{N}(-1, 1)$ , *i.e.*, the scores are Gaussian-distributed around +1 and -1, respectively. In the second simulation, we draw classifier scores from  $p_1(s)$  and  $p_0(s)$  obtained from a real classifier (the distributions, not shown, are highly non-Gaussian). Results for both simulations are shown in Fig. 13a. Each curve shows the AP for a classifier as the evaluation set size (x-axis) varies. Different curves correspond to different category frequencies (0.001 to 0.1). Given a *fixed* classifier, we observe that the curves are ordered top-to-bottom by higher-to-lower  $f_c$ . This ordering shows that for a fixed classifier, AP tends to be lower for rarer categories, which is an intuitive and well-known trend (see [14] §3.2). More surprising is the finding that within a curve, *AP is consistently higher when the evaluation subset size is smaller*. This pattern exists at all frequency levels in both simulations.

Now moving from the simulated classification problem to real object detection data, we show  $AP^{bb}@75$  of a trained detector for three categories evaluated on random COCO `val2017` subsets of various sizes in Fig. 13b. The toaster category is one of the two rarer categories in COCO while cats and dining table appear more frequently. In each case we observe similar trends as in the earlier simulations.

Most categories in COCO are well-sampled like the cat and dining table categories and their AP has already converged on the 5k `val2017` set. Therefore overall AP does not vary much on COCO when comparing `val2017` to `test2017` results.

<sup>6</sup>The simulation code will be available on the LVIS website.



(a) AP of simulated classifiers as a function of the evaluation set size and the fraction of positive examples  $f_c$  (the number below each data point indicates the number of positives at that point, the shaded region indicates the standard error when averaged over 300 trials). The left plot shows the behavior of a random Gaussian classifier; the right shows a classifier that mimics the empirical score distribution of a trained classifier. While smaller  $f_c$  leads to decreased AP, we also observe a consistent decrease in AP as the evaluation set size increases (until convergence).

Figure 13. **AP bias** as the size of the evaluation set is varied.

subset size	AP	AP <sub>r</sub>	AP <sub>c</sub>	AP <sub>f</sub>
5k	24.8 $\pm$ 0.51	11.5 $\pm$ 1.71	25.5 $\pm$ 0.86	30.1 $\pm$ 0.28
10k	22.1 $\pm$ 0.31	10.5 $\pm$ 0.66	22.9 $\pm$ 0.56	29.2 $\pm$ 0.20
15k	20.8 $\pm$ 0.23	10.0 $\pm$ 0.54	21.7 $\pm$ 0.37	28.9 $\pm$ 0.11
20k (full)	18.4	8.8	18.7	27.2

(a) Fixed Mask R-CNN model (Table 3b best + scale jitter) evaluated on different size subsets of `test` v0.5 (average over 30 random subsets). AP on the 5k subset is similar to AP on `val` v0.5. As we increase subset size we observe a systematic decrease in all AP metrics consistent with the simulated and observed bias described in the main text.

Table 4. **Results on LVIS `test` v0.5** for different size subsets of `test` and for three different baseline models.

LVIS, unlike COCO, is not artificially balanced and therefore it contains a large number of rare categories. Therefore we expect to see a change in AP when moving from the small `val` v0.5 set to the larger `test` v0.5 set. We see exactly the predicted effects in Table 4a where we report the results of a fixed Mask R-CNN model on various sized subsets of the `test` set from 5k to 20k images. The results on the 5k subset size are in line with results on the 5k image `val` set. As the number of evaluation images increases, we observe that AP systematically decreases.

Despite the bias between `val` and `test` (or more generally evaluation sets of different sizes), we expect the ranking of different models on the `val` set and `test` to remain constant under typical conditions as it is not obvious how one would exploit the bias. In Table 4b we compare three models with distinctly different AP on the `val` set to each other on the `test` test. Indeed for at least these three models the ranking on `val` transfers exactly to the ranking on `test`.

### C.3. Comparing Models

Finally, to gain a sense for the statistical differences between the three models in Table 4b, we applied three standard hypothesis tests (paired t-test, random permutation test, and percentile bootstrap) to the mean of the per-category AP differences between pairs of models. As an example, for the ResNet-101 and ResNeXt-101 based models (`test` AP of 20.0 vs. 20.5, respectively) the paired t-

model	eval. set	AP	AP <sub>r</sub>	AP <sub>c</sub>	AP <sub>f</sub>
ResNet-50	val	24.4	14.5	24.3	28.4
	test	18.4	8.8	18.7	27.2
ResNet-101	val	26.0	15.8	26.1	29.8
	test	20.0	9.4	21.0	28.7
ResNeXt-101-32 $\times$ 8d	val	27.1	15.6	27.5	31.4
	test	20.5	9.8	21.1	30.0

(b) We compare how AP transfers for three different models (Table 3c) from `val` v0.5 to `test` v0.5. All AP metrics decrease but the ranking of the models remains consistent across `val` and `test`.

test and random permutation test return  $p = 0.0490$  and  $p = 0.0486$ , respectively. The percentile bootstrapped produced a 95% confidence interval of [0.002, 1.015], which excludes 0 (barely). These tests agree with each other (*i.e.*, they reject, at a 5% level, the null hypothesis that the mean difference in per-category AP values is zero) and provide some intuition for the statistical significance that might arise from a 0.5% absolute difference in AP on `test`.

### C.4. Summary

In existing class-balanced detection datasets, researchers have grown accustomed to AP transferring nearly perfectly between small validation sets (*e.g.*, 5k images) and larger test sets (20k images). In this section we demonstrated that when a dataset has a larger class imbalance there are at least two factors that cause AP estimated on smaller evaluation sets to be biased compared to larger evaluation sets. Empirically, this bias leads to higher AP on `val` v0.5 than on `test` v0.5. While a small validation set was unavoidable for LVIS v0.5, based on this analysis we may extend the validation set to include more images in release v1.

**Acknowledgements.** We would like to thank Ilija Radosavovic, Amanpreet Singh, Alexander Kirillov, Judy Hoffman, and Tsung-Yi Lin for their help during creation of LVIS. We thank the COCO Committee for granting us permission to annotate the COCO test set and Amanpreet Singh for help in creating the LVIS website.

## References

- [1] Fred Attneave and Malcolm D Arnoult. The quantitative study of shape and pattern perception. *Psychological bulletin*, 1956. 7
- [2] Léon Bottou. Stochastic gradient descent tricks. In *Neural networks: Tricks of the trade*, pages 421–436. Springer, 2012. 10
- [3] Marc Brysbaert, Amy Beth Warriner, and Victor Kuperman. Concreteness ratings for 40 thousand generally known english word lemmas. *Behavior research methods*, 2014. 6
- [4] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The Cityscapes dataset for semantic urban scene understanding. In *CVPR*, 2016. 2
- [5] Piotr Dollár, Christian Wojek, Bernt Schiele, and Pietro Perona. Pedestrian detection: An evaluation of the state of the art. *TPAMI*, 2012. 2
- [6] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The PASCAL Visual Object Classes (VOC) Challenge. *IJCV*, 2010. 2
- [7] Li Fei-Fei, Rob Fergus, and Pietro Perona. One-shot learning of object categories. *TPAMI*, 2006. 2
- [8] Ross Girshick, Ilya Radosavovic, Georgia Gkioxari, Piotr Dollár, and Kaiming He. Detectron. <https://github.com/facebookresearch/detectron>, 2018. 8
- [9] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch sgd: Training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017. 9
- [10] Bharath Hariharan and Ross Girshick. Low-shot visual recognition by shrinking and hallucinating features. In *ICCV*, 2017. 8
- [11] Haibo He and Edwardo A Garcia. Learning from imbalanced data. *IEEE Transactions on knowledge and data engineering*, 2009. 10
- [12] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask R-CNN. In *ICCV*, 2017. 8, 9
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 2
- [14] Derek Hoiem, Yodsawalai Chodpathumwan, and Qieyun Dai. Diagnosing error in object detectors. In *ECCV*. 2012. 11
- [15] Sergey Ioffe and Christian Szegedy. Batch normalization: accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. 2
- [16] Alexander Kirillov, Kaiming He, Ross Girshick, Carsten Rother, and Piotr Dollár. Panoptic segmentation. In *CVPR*, 2019. 7
- [17] Alex Krizhevsky, Ilya Sutskever, and Geoff Hinton. ImageNet classification with deep convolutional neural networks. In *NIPS*, 2012. 2
- [18] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Malloci, Tom Duerig, et al. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *arXiv preprint arXiv:1811.00982*, 2018. 3
- [19] Yann LeCun, Bernhard Boser, John S Denker, Donnie Henderson, Richard E Howard, Wayne Hubbard, and Lawrence D Jackel. Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1989. 2
- [20] Yann LeCun, Corinna Cortes, and Christopher J.C. Burges. The MNIST database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998. 2
- [21] Marc Liberman. Reproducible research and the common task method. Simmons Foundation Lecture <https://www.simonsfoundation.org/lecture/reproducible-research-and-the-common-task-method/>, 2015. 2
- [22] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *CVPR*, 2017. 9
- [23] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft COCO: Common objects in context. In *ECCV*, 2014. 1, 2
- [24] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. COCO detection evaluation. <http://cocodataset.org/#detection-eval>, Accessed Oct 30, 2018. 2, 3
- [25] Dhruv Mahajan, Ross Girshick, Vignesh Ramanathan, Kaiming He, Manohar Paluri, Yixuan Li, Ashwin Bharambe, and Laurens van der Maaten. Exploring the limits of weakly supervised pre-training. In *ECCV*, 2018. 10
- [26] David Martin, Charless Fowlkes, Doron Tal, and Jitendra Malik. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. In *ICCV*, 2001. 2, 7
- [27] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013. 10
- [28] George Miller. *WordNet: An electronic lexical database*. MIT press, 1998. 4
- [29] Gerhard Neuhold, Tobias Ollmann, Samuel Rota Bulò, and Peter Kortscheder. The mapillary vistas dataset for semantic understanding of street scenes. In *ICCV*, 2017. 2
- [30] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *IJCV*, 2015. 2
- [31] Bryan C Russell, Antonio Torralba, Kevin P Murphy, and William T Freeman. Labelme: a database and web-based tool for image annotation. *IJCV*, 2008. 1
- [32] Li Shen, Zhouchen Lin, and Qingming Huang. Relay backpropagation for effective learning of deep convolutional neural networks. In *ECCV*, 2016. 10
- [33] Merrielle Spain and Pietro Perona. Measuring and predicting importance of objects in our visual world. Technical Report CNS-TR-2007-002, California Institute of Technology, 2007. 1
- [34] Grant Van Horn, Oisin Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, Pietro Perona, and Serge Belongie. The iNaturalist species classification and detection dataset. In *CVPR*, 2018. 2
- [35] Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. SUN database: Large-scale scene recognition from abbey to zoo. In *CVPR*, 2010. 1
- [36] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *CVPR*, 2017. 10
- [37] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Semantic understanding of scenes through the ADE20K dataset. *IJCV*, 2019. 2
- [38] George Kingsley Zipf. *The psycho-biology of language: An introduction to dynamic philology*. Routledge, 2013. 1