# MOSAICOS: A Simple and Effective Use of Object-Centric Images for Long-Tailed Object Detection

Cheng Zhang[1*]   Tai-Yu Pan[1*]   Yandong Li[2]   Hexiang Hu[3]
Dong Xuan[1]   Soravit Changpinyo[2]   Boqing Gong[2]   Wei-Lun Chao[1]

[1]The Ohio State University   [2]Google Research   [3]University of Southern California

## Abstract

*Many objects do not appear frequently enough in complex scenes* (e.g., *certain handbags in living rooms) for training an accurate object detector, but are often found frequently by themselves* (e.g., *in product images). Yet, these* object-centric *images are not effectively leveraged for improving object detection in* scene-centric *images. In this paper, we propose **Mosaic** of **O**bject-centric images as **S**cene-centric images (MOSAICOS), a simple and novel framework that is surprisingly effective at tackling the challenges of long-tailed object detection. Keys to our approach are three-fold: (i) pseudo scene-centric image construction from object-centric images for mitigating domain differences, (ii) high-quality bounding box imputation using the object-centric images' class labels, and (iii) a multi-stage training procedure. On LVIS object detection (and instance segmentation), MOSAICOS leads to a massive 60% (and 23%) relative improvement in average precision for rare object categories. We also show that our framework can be compatibly used with other existing approaches to achieve even further gains. Our pre-trained models are publicly available at* https://github.com/czhang0528/MosaicOS/.

## 1. Introduction

Detecting objects in complex daily scenes is a long-standing task in computer vision [17, 21, 54, 74]. With rapid advances in deep neural networks [30, 35, 41, 62, 64] and the emergence of large-scale datasets [43, 50, 61, 71, 85], there has been remarkable progress in detecting *common* objects (*e.g.*, cars, humans, *etc.*) [4, 29, 48, 49, 51, 56, 59, 91]. However, detecting *rare* objects (*e.g.*, unicycles, bird feeders, *etc.*) proves much more challenging due to the inherent limitation of training data. In particular, complex scenes in which an object appears pose another variation
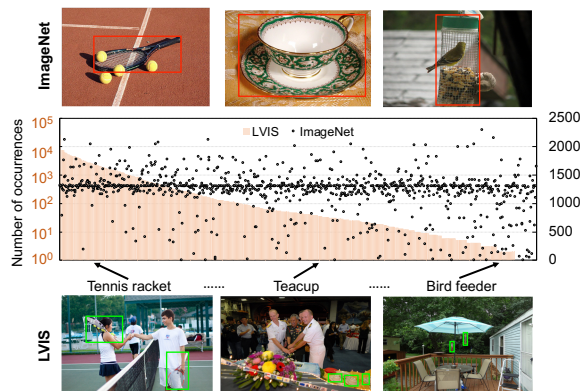


Figure 1. **Object frequencies in scene-centric and object-centric images.** Orange bars show the number of instances per class in the scene-centric LVIS v0.5 dataset [24]. Class indices are sorted by the instance numbers. Black dots show the number of images in the object-centric ImageNet datasets [13] for each corresponding class. *The two types of images have very different trends of object frequencies.* We also show three examples of both datasets, corresponding to frequent, common, and rare classes in LVIS. **Red** and **green** boxes indicate the objects. *These two types of images have different focuses and object sizes.*

factor that is too diverse to capture from a small amount of data [24, 73, 90]. Developing algorithms to overcome such a "long-tailed" distribution of object instances in *scene-centric* images (SCI) [24, 43, 50, 61] has thus attracted a flurry of research interests [47, 66, 77].

Fortunately, the uncommon objects in scene-centric images often appear more frequently in *object-centric* images (OCI) in which the objects of interest occupy the main and most salient regions (*e.g.*, product images). For example, given "bird feeder" as query, a popular image search engine (*e.g.*, Google Images) mostly retrieves object-centric "bird feeder" results. Similarly, curated object recognition datasets such as ImageNet [13] contain more than a thousand object-centric "bird feeder" images, nearly a hundred times more than scene-centric images from LVIS v0.5 [24]. We further illustrate this point in Figure 1, in which a discrepancy in frequencies of the same objects in ImageNet

and LVIS can generally be observed (see § 3 for details).

Can we leverage such abundant object-centric images to improve long-tailed object detection? The most common approach to this is to leverage these images for pre-training the object detector's backbone [28, 29, 91]. While this general approach may benefit various tasks beyond object detection, it is highly data-intensive and does not take care of the domain gap between the pre-training and downstream tasks (see § 6.3 for analysis). As a result, they do not always improve the object detection accuracy [28, 91].

In this paper, we propose MOSAICOS (**Mosaic** of **O**bject-centric images as **S**cene-centric images), a simple and effective framework to leverage object-centric images for object detection. MOSAICOS directly uses object-centric images during the *training* of object detectors. There are three key ingredients. The first one is the construction of *pseudo* scene-centric images from object-centric images using mosaic[1]. The second one is the imputation of bounding box annotations using image class labels. The third ingredient is a multi-stage training procedure for learning from both gold scene-centric and synthesized pseudo scene-centric annotations. Figure 2 illustrates our framework.

Our use of mosaic and bounding box imputation to construct *pseudo* scene-centric images from object-centric images tackles two key challenges in leveraging object-centric images for object detection. The first challenge is a "domain" gap between object-centric and scene-centric images: an object-centric image usually contains fewer (but bigger) object instances and a less complex background and this discrepancy is believed to unfavorably hinder knowledge transfer [78, 86]. The second challenge is missing detection annotations: object-centric images, either from the Internet or object recognition datasets (*e.g.*, ImageNet), are not *perfectly object-centric*, usually provided *without* accurate object localization in the form of bounding boxes.

Our proposed framework leads to significant accuracy gains in long-tailed object detection and instance segmentation on LVIS [24], using object-centric images from ImageNet [13] and the Internet. In particular, for the task of object detection for rare objects, we observe a significant boost from 13% to over 20% in average precision. For the task of instance segmentation, our approach even improves the accuracy on common objects. More importantly, unlike the baseline approaches, we do so without sacrificing the accuracy on the frequent classes. Finally, we also explore combining our approach with existing techniques [58] that results in even better performance.

**Our main contributions** are summarized as follows:
- Bringing the best of object-centric images to the long-tailed object detection on scene-centric images.

- Algorithms for mosaicking and pseudo-labeling to mitigate the domain discrepancy between two image types.
- A multi-stage training framework that leverages the pseudo scene-centric images (from object-centric images) to improve the detector on scene-centric images.
- Extensive evaluation and analysis of the proposed approach on the challenging LVIS benchmark [24].

## 2. Related Work

**Long-tailed object detection** has attracted increasing attentions recently. The challenge is the drastically low accuracy for detecting rare objects. Most existing works develop training strategies or objectives to address this [24, 34, 47, 55, 58, 66, 68, 76, 77, 80]. Wang et al. [76] found that the major performance drop is by mis-classification, suggesting the applicability of class-imbalanced classification methods (*e.g.*, re-weighting, re-sampling) [6–8, 11, 11, 23, 26, 26, 38, 40, 50, 71, 73, 88]. Different from them, we study an alternative and orthogonal solution to the problem (*i.e.*, exploiting abundant object-centric images).

**Weakly-supervised or semi-supervised object detection** learns or improves object detectors using images with weak supervision (*e.g.*, image-level labels) [3, 14, 19, 45, 72] or even without supervision [19, 37, 46, 63, 91]. They either leverage scene-centric images or detect only a small number of common classes (*e.g.*, classes in Pascal VOC [15], MSCOCO [50], or ILSVRC [60]). Our work can be seen as weakly supervised object detection, but we focus on the challenging long-tailed detection with more than $1,000$ objects. Meanwhile, we leverage object-centric images, which is different from scene-centric images in both appearances and layouts. The most related work is [55], which uses the YFCC-100M dataset [71] (Flickr images) to improve the detection on LVIS [24]. However, YFCC-100M contains both object-centric images and scene-centric images and a non-negligible label noises. Thus, [55] employs more sophisticated data pre-processing and pseudo-labeling steps, yet our approach achieves higher accuracy (see Table 5).

Other works use object-centric images to expand the label space of the object detector [31–33, 42, 57, 69, 70]. Such approaches mostly only use object-centric images to learn the last fully-connected classification layer, instead of improving the features extractor. In contrast, our approach can improve the feature extractor, and successfully transfer knowledge to long-tailed instance segmentation.

## 3. Scene-Centric vs. Object-Centric Images

Images taken by humans can roughly be categorized into object-centric and scene-centric images. The former captures objects of interest (*e.g.*, cats) and usually contains just one salient class whose name is used as the image label. The later captures a scene and usually contains multiple object

---

[1]Mosaic was exploited in [4, 10, 89], but mainly to combine multiple *scene-centric* images to simulate smaller object sizes or to increase the scene complexity, not to turn object-centric images into scene-centric ones.

instances of different classes in a complex background.

Recent object detection methods mainly focus on scene-centric images [24, 50, 61]. Since scene-centric images are not intended to capture specific objects, *object frequencies in our daily lives will likely be reflected in the images.* As such, the learned detector will have a hard time detecting rare objects: it just has not seen sufficient instances to understand the objects' appearances, shapes, variations, etc.

In contrast, humans tend to take object-centric pictures that capture interesting (and likely uncommon, rare) objects, especially during events or activities (*e.g.*, bird watching, *etc.*). Thus, a rare object in our daily lives may occur more often in the online object-centric images.

**Discrepancy w.r.t. object frequencies.** We compare object frequencies of the ImageNet [13] and LVIS (v0.5) [24] datasets. The former retrieved images from the Internet by querying search engines using the object class names (thus object-centric). Whereas the later used MSCOCO [50] dataset, which collects daily scene images with many common objects in a natural context (thus scene-centric).

The full ImageNet has $21,841$ classes, whereas LVIS has around $1,230$ classes. Using the WordNet synsets [53], we can match $1,025$ classes ($1,016$ classes are downloadable) between them. Figure 1 shows the number of object instances per class in LVIS and the number of images per corresponding class in ImageNet. It presents *a huge difference between object frequencies of these two datasets.* For example, ImageNet has a balanced distribution across classes and LVIS is extremely long-tailed. Even for rare classes in LVIS (those with $< 10$ training images), ImageNet usually contains more than $1,000$ images. Such a difference offers an opportunity to resolve the long-tailed object detection in scene-centric images via the help of object-centric images.

**Discrepancy w.r.t. visual appearances and contents.** Beyond frequencies, these two types of images also have other, less favorable discrepancies. The obvious one is the number of object instances per image. LVIS on average has **12.1** labeled object instances per image (the median number of instances per image is **6**). While most of the ImageNet images are not annotated with object bounding boxes, according to a subset of images used in the ILSVRC detection challenge [60], each image has **2.8** object instances. The larger number of object instances, together with the intention behinds the images, implies that scene-centric images also have *smaller* objects in size and more complex backgrounds. This type of discrepancies, contrast to that in object frequencies, is not favorable for leveraging the object-centric images, and may lead to negative transfer [78, 86].

## 4. Overall Framework

To better leverage object-centric images for object detection, we present a novel learning framework, which includes three simple[2] **yet effective components** to handle

(a) the domain gap between two image sources;
(b) the missing bounding box labels;
(c) the integration of both image sources for training.

Figure 2 gives an illustration of our framework. Concretely, the framework begins with pre-training an object detector using the accurately labeled scene-centric images. Any object detector can be applied. Without loss of generality, we focus on Faster R-CNN [59], one of the most popular object detectors in the literature. The pre-trained object detector serves for two purposes: it can help impute the missing boxes in object-centric images; it will be used as the initialization for training with the object-centric images.

To turn object-centric images into training examples for object detection, we must handle both (a) and (b). *We postpone the details of these two components to § 5.* For now, let us assume that we have processed and labeled object-centric images with pseudo ground-truth boxes and class labels like the labeled scene-centric images. To differentiate from the original object-centric images, we call the new images *pseudo scene-centric images* (see Figure 2).

The pseudo scene-centric images may still have domain gaps from real scene-centric images, to which the learned detector will finally be applied. Besides, the pseudo ground-truth boxes may contain noises (*e.g.*, wrong locations). To effectively learn from these images (especially for rare objects) while not sacrificing the detector's ultimate accuracy on identifying and locating objects, we propose to *fine-tune the pre-trained object detector via two stages.* In what follows, we first give a brief review of object detection.

**Backgrounds on object detection.** An object detector has to identify objects with their class names and locate each of them by a bounding box. Taking Faster R-CNN [59] as an example, it first generates a set of object proposals (usually around $512$) that likely contain objects of interest. This is done by the region proposal network (RPN) [59]. Faster R-CNN then goes through each proposal to identify its class (can be "background") and refine the box location and size.

The entire Faster R-CNN is learned with three loss terms

$$\mathcal{L} = \mathcal{L}_{\text{rpn}} + \mathcal{L}_{\text{cls}} + \mathcal{L}_{\text{reg}}, \tag{1}$$

where $\mathcal{L}_{\text{rpn}}$ is for RPN training, $\mathcal{L}_{\text{cls}}$ is for multi-class classification, and $\mathcal{L}_{\text{reg}}$ is for box refinement.

**Two-stage fine-tuning.** Given the pre-trained detector, we first fine-tune it using the pseudo scene-centric images that are generated from object-centric images (see § 5). We then fine-tune it using the labeled scene-centric images. We separate the two image sources since they are still different

---

[2]We claim our approach to be "simple" as it employs simple methods to address the fundamental challenges. Pseudo-labeling is an essential step to use weakly-supervised data, and we apply simple fixed locations. We apply mosaicking and multi-stage training to bridge the domain gap, instead of applying sophisticated methods like domain adversarial training [18].
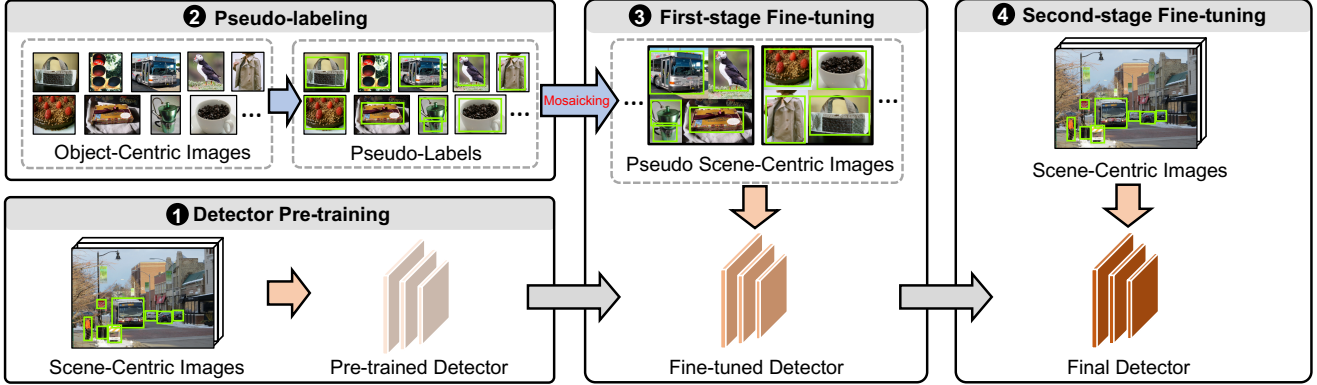
Figure 2. **Our MOSAICOS framework for leveraging object-centric images for long-tailed object detection.** It consists of four stages. ❶ **Detector Pretraining**: we pre-train an object detector using scene-centric images with gold-labeled box annotations. ❷ **Pseudo-labeling**: we construct pseudo scene-centric images from object-centric images using box annotation imputation (possibly using the pre-trained detector in stage 1) as well as mosaicking (stitching multiple images together). ❸ **First-stage Fine-tuning**: we fine-tune the pre-trained detector from stage 1 with pseudo scene-centric images from stage 2. ❹ **Second-stage Fine-tuning**: we further fine-tune the object detector from stage 3 using scene-centric images with gold-labeled box annotations again, similar to stage 1. Orange arrows indicate data feeding for training. Gray arrows indicate model cloning. Green boxes indicate the (pseudo & gold-labeled) box annotations.

in appearances and label qualities. The second stage helps adapt the detector back to real scene-centric images.

In both stages, all the three loss terms in Equation 1 are optimized. We do not freeze any parameters except the batch-norm layers [36] in the backbone feature network which are kept frozen by default. We will compare to single-stage fine-tuning with both images and fine-tuning using only $\mathcal{L}_{cls}$ for pseudo scene-centric images in § 6.

## 5. Creating Pseudo Scene-Centric Images

We now focus on the missing components of our framework: generating pseudo scene-centric images from object-centric images. Our goal is to create images that are more *scene-centric-like* and label them with pseudo ground truths. We collect object-centric images from two sources: ImageNet [13] and Google Images. See § 6.1 for details.

### 5.1. Assigning Pseudo-Labels

Each object-centric image has one object class label, but no bounding box annotations. Some images may contain multiple object instances and classes, in which the class label only indicates the most salient object. Our goal here is to create a set of pseudo ground-truth bounding boxes that likely contain the object instances for each image, and assign each of them a class label, such that we can use the image to directly fine-tune an object detector.

There are indeed many works on doing so, especially those for weakly-supervised and semi-supervised object detection [3, 14, 19, 45, 46, 55, 72]. The purpose of this subsection is therefore not to propose a new way to compare with them, but to investigate approaches that are more effective and efficient in a large-scale long-tailed setting. Specifically, we investigate five methods that do not require an

extra detector or proposal network beyond the pre-trained one. *As will be seen in § 6.2, imputing the box class labels using the image class label is the key to success.* Figure 3 illustrates the difference of these methods. Please see the supplementary material for other possibilities.

**Fixed locations (F).** We simply assign some fixed locations of an object-centric image to be the pseudo ground-truth boxes, regardless of the image content. The hypothesis is that many of the object-centric images may just focus on one object instance whose location is likely in the centre of the image (*i.e.*, just like those in [16, 22]). Specifically, we investigate the combination of the whole image, the center crop, and the four corner crops: in total **six** boxes per image. The height and width of the crops are $80\%$ of the original image. We assign each box the image class label.

**Trust the pre-trained detector (D).** We apply the pre-trained detector learned with the scene-centric images to the object-centric images, and treat the detected boxes and predicted class labels as the pseudo-labels. Specifically, we keep all the detection of confidence scores $> 0.5$[3]. We apply non-maximum suppression (NMS) among the detected boxes of each class using an IoU (intersection-over-union) threshold $0.5$. By doing so, every image will have boxes of different sizes and locations, labeled with different classes.

**Trust the pre-trained detector & image class labels (D†).** One drawback of the above method is its tendency to assign high-frequency labels, a notorious problem in class-imbalanced learning [38, 66, 84]. For instance, if "bird" is a frequent class and "eagle" is a rare class, the detector may correctly locate an "eagle" in the image but assign the label "bird" to it. To resolve this issue, we choose to trust the box locations generated by the above method but assign

---

[3]$0.5$ is the default threshold for visualizing the detection results.
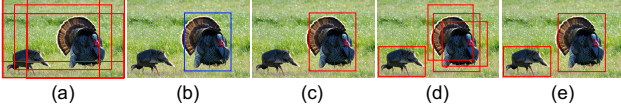
4

Figure 3. **A comparison of pseudo-label generation:** (a) fixed locations, (b) trust the detector, (c) trust the detector + image labels, (d) trust the calibrated detector + image labels, and (e) localization by region removal. The image label is "turkey," a rare class in LVIS. Red/blue boxes are labeled as "turkey"/other classes.

each box the image class label instead of the predicted class labels. In other words, a box initially labeled as "bird" is replaced by the label "eagle" if "eagle" is the image label. The rationale is that in an object-centric image, most of the object instances belong to the image's class.

**Trust the calibrated detector & image class labels ($D\ddagger$).** Another way to resolve the above issue is to set for each class a different confidence threshold[4]. The rational is that a classifier trained with long-tailed data tends to assign lower probabilities to minor classes of scarce training data [6, 38, 39, 84]. We thus reduce the threshold for each class according to its number of training images. Let $N_{max}$ be the size of the most major class and let $N_c$ be the size of class $c$, we apply a threshold $0.5 \times (N_c/N_{max})^\gamma$ for class $c$, inspired by [52]. We set $\gamma = 0.5$ according to validation. Compared to the vanilla "trust the detector," this method will generate more boxes for common and rare classes. We again replace their detected labels by the image class label[5].

**Localization by region removal (LORE).** We investigate yet another way for pseudo-labeling, taking the following intuition: an image classifier should fail to predict the correct label if the true object regions are removed. To this end, we first train a ResNet-50 [27] image classifier using our object-centric image pool[6]. We then collect the pretrained detector's predicted boxes on these images, trusting the box locations but not the class labels. We sort these boxes by how much removing each boxed region alone reduces the image classifier's confidence on the image label. We then remove these boxed regions *in turn* until the classifier fails to predict the true label. The bounding boxes of the removed regions are then collected as the pseudo ground truths for the image. We assign each box the image class label. Please see the supplementary material for details.

### 5.2. Synthesizing Pseudo Scene-centric Images

We apply a simple technique, *i.e.*, *image mosaic*, to make object-centric images more scene-centric, in terms of appearances, contents, and layouts. Concretely, we stitch mul-

tiple object-centric images together to obtain a new image that contains more object instances and classes, smaller object sizes, and more complex background. Specifically, we stitch $2 \times 2$ images together, which are sampled either randomly within a class, or randomly from the entire image pool. We do not apply sophisticated stitching tools like [5, 25, 87] but simply concatenate these images one-by-one. The resulting images are thus more like *mosaics*, having artifacts along the stitched boundaries (see Figure 2).

## 6. Experiments

We conduct experiments and analysis for MOSAICOS, on the tasks of long-tailed Object Detection (OD) and Instance Segmentation (IS). We begin by introducing the experimental setup (§ 6.1), then present the main object detection results as well as detailed ablation studies (§ 6.2, § 6.3), and finally show additional results that evaluate our model on instance segmentation and the other dataset (§ 6.4). *We include qualitative results in the supplementary material.*

### 6.1. Setup

**Long-tailed OD & IS datasets and metrics.** We evaluate our approach on LVIS instance segmentation benchmark [24]. We focus on v0.5 as most existing works, and report additional key results on v1.0 (more in the supplementary). LVIS v0.5 contains $1,230$ entry-level object categories with around 2 million high-quality annotations. The training set contains all the classes with a total of $57,623$ image; the validation set contains 830 classes with a total of $5,000$ images. The categories are naturally long-tailed distributed and are divided into three groups based on the number of training images per class: rare (1-10 images), common (11-100 images), and frequent (>100 images). *All results are reported on the validation set.* We adopt the standard mean average precision (AP) metric in LVIS [24]. *We specifically focus on object detection using the standard bounding box evaluation, $AP^b$.* The AP on the rare, common, and frequent classes ($AP_r^b$, $AP_c^b$, $AP_f^b$) are also reported separately.

**Object-centric data sources.** We mainly use images from two sources: ImageNet [13] and Google Images [2]. ImageNet is a classification benchmark. Most people use its $1,000$ categories version in ILSVRC [60] and treat it as the standard dataset for backbone pre-training in various computer vision tasks. The full version of ImageNet has $21,842$ classes. In LVIS and ImageNet, each category has a unique WordNet synset ID [53], and we are able to match $1,016$ LVIS classes and retrieve the corresponding images from ImageNet (in total, $769,238$ images). Besides, we retrieve images via Google by querying with class names and descriptions provided by LVIS. Such a text-to-image search returns hundreds of iconic images and we take the top 100 for each of the $1,230$ classes.

---

[4]For Faster R-CNN, each RPN proposal can lead to multiple detected boxes, one for each class whose probability is larger than the threshold.

[5]Without doing so, the approach can hardly improve "trust the detector" due to more noisy boxes being included. See the supplementary for details.

[6]That is, we train the classifier with these images, and then apply this classifier back to these images (after some regions are removed).

5

**Implementation.** We use Faster R-CNN [59][7] as our base detector and ResNet-50 [30] with a Feature Pyramid Network (FPN) [48] as the backbone, which is pre-trained on ImageNet (ILSVRC) [60]. Our base detector is trained on the LVIS training set with *repeated factor sampling*, following the standard training procedure in [24] (1x schedule). To fairly compare with our fine-tuning results, we further extend the training process with another 90K iterations and select the checkpoint with the best $AP^b$ as **Faster R-CNN⋆**. The following experiments are initialized by Faster R-CNN⋆. See the supplementary for more details.

**Baselines.** We compare to the following baselines:

- **Self-training** is a strong baseline for semi-supervised learning [44]. We follow the state-of-the-art self-training method for detection [91] and use Faster R-CNN⋆ to create pseudo-labels on the object-centric images, same as "trust the pre-trained detector". We then fine-tune Faster R-CNN⋆ using both pseudo scene-centric and LVIS images for 90K iterations, with the normalized loss [91].
- **Single-stage fine-tuning** fine-tunes Faster R-CNN⋆ with both pseudo scene-centric and LVIS images in one stage. In each mini-batch, we have 50% of data from each source. Different ratios do not lead to notable differences.
- **DLWL** [55] is the state-of-the-art method that uses extra unlabeled images from the YFCC-100M [71].

For self-training and single-stage training, we perform $2 \times 2$ mosaicking to create pseudo scene-centric images.
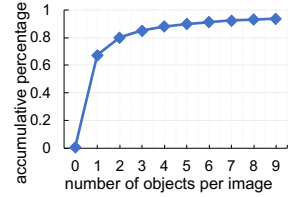
**Variants of** MOSAICOS. (a) We compare different object-centric image sources and their combinations. (b) We compare with or without mosaicking. (c) For mosaicking, we compare stitching images from the same classes (so the artifacts can be reduced) or from randomly selected images. (d) We compare different ways to generate pseudo-labels (see § 5.1). (e) We study fine-tuning with pseudo scene-centric images using only the classification loss $\mathcal{L}_{cls}$.

## 6.2. Results on Object Detection

**Main results.** Table 1 summarizes the results. It shows that the model trained with pseudo-labels generated by six fixed location (F) has a very competitive performance comparing to other strategies. We therefore consider it as the default pseudo-labeling method given its simplicity and effectiveness. Meanwhile, our two-stage approach with object-centric images outperforms Faster R-CNN⋆ (and Faster R-CNN) notably. On $AP^b_r$ for rare classes, our best result of 20.25% is $\sim 7.2\%$ higher than Faster R-CNN⋆, justifying our motivation: object-centric images that are resistant to the long-tailed nature of object frequencies can improve object detection in scene-centric images.

**Mosaicking is useful** (red in Table 1). A simple $2 \times 2$

---

---



Figure 4. **# of objects per object-centric image found by LORE.** We use ImageNet images. Y-axis is the accumulative percentage of images whose objects are no more than the X-axis number.

stitching leads to a notable gain: $\sim 1.8\%$ at $AP^b_r$, supporting our claim that making object-centric images similar to scene-centric images is important. Indeed, according to § 3, a $2 \times 2$ stitched image will have around 12 objects, very close to that in LVIS images. Stitching images from different classes further leads to a $\sim 1.0\%$ gain (green in Table 1).

**Fixed-location boxes are effective** (blue in Table 1). By comparing different ways for pseudo-labels, we found that both localization by region removal (L) and the simple six fixed locations (F) lead to strong results without querying the pre-trained detector. Using six fixed locations slightly outperforms using one location (S) (*i.e.*, the image boundary), probably due to the effect of data augmentation. All the three methods significantly outperform "trust the pre-trained detector" (D), and we attribute this to the poor pre-trained detector's accuracy on rare classes: it either cannot identify rare classes from object-centric images or is biased to detect frequent classes. **By replacing the detected classes with the image labels** and/or further calibrating the detector for more detected boxes, *i.e.*, "trust the (calibrated) pre-trained detector and image label" (D†, D‡), we see a notable boost, which supports our claim. Nevertheless, they are either on par with or worse than fixed locations (F), especially on $AP^b_r$ for rare classes, again showing the surprising efficacy of the simple method. We note that, both D† and D‡ are specifically designed in this work for long-tailed problems and should not be seen as existing baselines.

To further analyze why fixed locations work well, we check the numbers of boxes LORE finds per image. LORE keeps removing regions until the classifier fails to classify the image. The number of regions it found is thus an estimation of the number of target objects (those of the image label) in the image. Figure 4 shows the accumulative number of images whose object numbers are no more than a threshold: $\sim 70\%$ of ImageNet images have one target object instance, suggesting that it may not be necessary to locate and separate object instances in pseudo-labeling.

**Self-training and single-stage fine-tuning.** As shown in Table 1, self-training (with D and loss normalization [91]) outperforms Faster R-CNN⋆ on $AP^b_r$. As self-training is sensitive to the pseudo-label quality, we also apply the fixed locations (F) to it and achieve improvement. By comparing it to single-stage fine-tuning (with F), we see the benefit of loss normalization between the two image sources.

By comparing self-training to its counterparts in two-stage fine-tuning (with D and F), we however find that two-

Table 1. **Comparison of object detection on LVIS v0.5 validation set. OCIs:** object-centric images sources (IN: ImageNet, G: Google). **Mosaic:** ✓means 2×2 image mosaicking. **Hybrid:** ✓means stitching images from different classes. **P-GT:** ways to generate pseudo-labels (**F**: six fixed locations, **D**: trust the detector, **D†**: trust the detector and image label, **D‡**: trust the calibrated detector and image class label, **L** (LORE): localization by region removal, and **S**: a single box of the whole image). The best result per column is in bold.

| | OCIs | Mosaic | Hybrid | P-GT | $AP^b$ | $AP^b_{50}$ | $AP^b_{75}$ | $AP^b_r$ | $AP^b_c$ | $AP^b_f$ |
|---|---|---|---|---|---|---|---|---|---|---|
| Faster R-CNN | - | - | - | - | 23.17 | 38.94 | 24.06 | 12.64 | 22.40 | 28.33 |
| Faster R-CNN⋆ | - | - | - | - | 23.35 | 39.15 | 24.15 | 12.98 | 22.60 | 28.42 |
| Self-training [91] | IN | ✓ | ✓ | D | 22.71 | 38.22 | 23.79 | 14.52 | 21.41 | 27.61 |
| | IN | ✓ | ✓ | F | 23.46 | 39.03 | 24.82 | 16.20 | 22.19 | 27.94 |
| Single-stage | IN | ✓ | ✓ | F | 20.09 | 35.34 | 20.27 | 12.96 | 19.08 | 24.20 |
| MOSAICOS (Two-stage) | IN | ✗ | ✗ | F | 24.27 | 40.30 | 25.61 | 16.97 | 23.29 | 28.42 |
| | IN | ✓ | ✗ | F | 24.48 | 40.12 | 25.65 | 18.76 | 23.26 | 28.29 |
| | IN | ✓ | ✓ | D | 23.04 | 38.97 | 23.72 | 13.93 | 21.51 | 28.14 |
| | IN | ✓ | ✓ | D† | 24.66 | 40.31 | 25.99 | 17.45 | 23.62 | 28.83 |
| | IN | ✓ | ✓ | D‡ | 24.93 | 40.48 | 26.71 | 19.31 | 23.51 | **28.95** |
| | IN | ✓ | ✓ | S | 24.59 | 40.20 | 25.78 | 19.13 | 23.35 | 28.32 |
| | IN | ✓ | ✓ | L | 24.83 | 40.58 | 26.27 | 20.06 | 23.25 | 28.71 |
| | IN | ✓ | ✓ | F | 24.75 | 40.44 | 26.09 | 19.73 | 23.44 | 28.39 |
| | IN+G | ✓ | ✓ | F | **25.01** | **40.76** | **26.46** | **20.25** | **23.89** | 28.32 |

Table 2. **Losses in the first fine-tuning stage.**

| Losses | $AP^b$ | $AP^b_r$ | $AP^b_c$ | $AP^b_f$ |
|---|---|---|---|---|
| $\mathcal{L}_{cls}$ | 24.53 | 18.87 | 23.07 | 28.61 |
| $\mathcal{L}_{rpn} + \mathcal{L}_{cls} + \mathcal{L}_{reg}$ | 24.75 | 19.73 | 23.44 | 28.39 |

Table 3. **Object detection on LVIS v0.5.** We use ImageNet + Google Images. MSCOCO: for pre-training. [58]: balanced loss.

| | MSCOCO | [58] | $AP^b$ | $AP^b_r$ | $AP^b_c$ | $AP^b_f$ |
|---|---|---|---|---|---|---|
| BaGS [47] | ✓ | | 25.96 | 17.65 | 25.75 | 29.54 |
| TFA [77] | | | 24.40 | 16.90 | 24.30 | 27.70 |
| MOSAICOS | | | 25.01 | 20.25 | 23.89 | 28.32 |
| | ✓ | | 26.28 | 17.37 | 26.13 | 30.02 |
| | | ✓ | 26.83 | **21.00** | 26.31 | 29.81 |
| | ✓ | ✓ | **28.06** | 19.11 | **28.23** | **31.41** |

Table 4. **The importance of mosaicking object-centric images.** SCI: object-centric images in the original LVIS training set.

| | $AP^b$ | $AP^b_r$ | $AP^b_c$ | $AP^b_f$ |
|---|---|---|---|---|
| Faster R-CNN⋆ | 23.35 | 12.98 | 22.60 | 28.42 |
| Stitching SCI [4] | 23.83 | 13.99 | 23.02 | **28.76** |
| Stitching SCI [10] | 23.58 | 14.00 | 22.58 | 28.66 |
| Stitching cropped SCI [89] | 23.55 | 13.40 | 23.04 | 28.26 |
| MOSAICOS | **24.75** | **19.73** | **23.44** | 28.39 |

stage fine-tuning leads to higher accuracy in most cases. This demonstrates the benefit of separating image sources in fine-tuning, in which the second stage adapts the detector back to accurately labeled true scene-centric images.

**The amount of object-centric data** (brown in Table 1). As ImageNet only covers $1,016$ classes of LVIS, we augment it with $100$ Google images per class for all the $1,230$ LVIS classes. We see another $0.5\%$ gain at rare classes ($AP^b_r$).

*For the following analyses besides Table 1, we will focus on our approach with two-stage fine-tuning, ImageNet object-centric images, $2 \times 2$ mosaic with images from multiple classes, and fixed locations (F) as the pseudo-labels.*

**Losses in fine-tuning.** We compare using all three losses of Faster R-CNN (*i.e.*, $\mathcal{L}_{rpn} + \mathcal{L}_{cls} + \mathcal{L}_{reg}$) or just the classification loss (*i.e.*, $\mathcal{L}_{cls}$) in the first-stage fine-tuning with pseudo scene-centric images. Table 2 shows the results. The former outperforms the latter on three out of four metrics. This tells that, while the pseudo boxes do not accurately bound the objects, learning the RPN and box refinement with them (*e.g.*, to predict a high objectness score) is still beneficial.

**Other baselines.** We compare to state-of-the-art methods that use no extra object-centric images in Table 3. We obtain comparable or better results, especially on rare classes.

**Compatibility with existing efforts.** Our approach is compatible with recent efforts in better backbone pre-training [47] and advanced training objectives (*e.g.*, [58]). For instance, following BaGS [47] to pre-train the backbone using MSCOCO images, we achieve an improved 26.28 $AP^b$ (see Table 3). Further incorporating the balanced loss [58] into the second-stage fine-tuning boosts $AP^b$ to 28.06.

## 6.3. Detailed Analysis of MOSAICOS

**The importance of object-centric images.** In Table 1, we show that even without mosaic, the use of object-centric images already notably improves the baseline ($AP^b$: 24.27 vs. 23.35). We further investigate the importance of mosaic of object-centric images: our use of mosaic is different from [4, 10, 89], which stitch scene-centric images in the training set to simulate smaller objects or increase the scene complexity. We apply their methods to stitch LVIS images and study two variants: stitching scene-centric images [4, 10] or the cropped objects [89] from them. Table 4 shows that MOSAICOS surpasses both variants on rare and common objects, justifying the importance of incorporating ample object-centric images to capture the diverse appearances of objects, especially for rare objects in scene-centric images.

**Does the quality of data sources matter?** DLWL [55] uses YFCC-100M [71], a much larger data source than Im-

Figure 5. **A comparison of object-centric image sources.** We show images of a rare class ("ax") in LVIS. ImageNet [13] (right) and Google Images [2] (middle) give images with more salient "ax" inside, while Flickr [1] (left) gives more noisy images, either with very small axes, cluttered backgrounds, or even no axes.

Table 5. **Comparison of object detection on LVIS v0.5 using different extra data sources.** G: Google Images. IN: ImageNet.

|  | Data | $AP^b$ | $AP^b_r$ | $AP^b_c$ | $AP^b_f$ |
|---|---|---|---|---|---|
| Faster R-CNN⋆ | – | 23.35 | 12.98 | 22.60 | 28.42 |
| DLWL [55] | YFCC-100M | 22.14 | 14.21 | - | - |
| MOSAICOS | Flickr | 24.05 | 16.17 | 23.06 | 28.43 |
|  | G | 24.45 | 19.09 | 23.27 | 28.08 |
|  | IN | 24.75 | 19.73 | 23.44 | 28.39 |

Table 6. **Object detection on the 176 overlapped classes** between ImageNet-1K (ILSVRC) and LVIS v0.5.

|  | $AP^b$ | $AP^b_r$ | $AP^b_c$ | $AP^b_f$ |
|---|---|---|---|---|
| # Category | 176 | 21 | 84 | 71 |
| Faster R-CNN⋆ | 26.05 | 14.78 | 23.92 | 31.16 |
| MOSAICOS | **27.50** | **21.16** | **25.45** | **31.80** |

ageNet. YFCC-100M images are mainly collected from Flickr, which mixes object-centric and scene-centric images and contains higher label noises. DLWL [55] thus develops sophisticated pre-processing and pseudo-labeling steps. In contrast, we specifically leverage *object-centric* images that have higher object frequencies and usually contain only single object classes (the image labels), leading to a much simpler approach. As shown in Table 5, our method (with IN) outperforms DLWL by a large margin: $> 5.5\%$ at $AP^b_r$. We attribute this to our ways of strategically collecting object-centric images and stitching them to make them scene-centric-like. The fact that we identify a better data source should not lead to an impression that we merely solve a simpler problem, but an evidence that selecting the right data source is crucial to simplify a problem. Figure 5 illustrates the difference among these sources.

For a fair comparison to [55] in terms of the algorithms, we also investigate Flickr images. Since [55] does not provide their processed data, we directly crawl Flickr images (100 per class) and re-train our algorithm. We achieve $24.05/16.17$ $AP^b/AP^b_r$, better than DLWL. Using pure Google images beyond ImageNet can achieve $24.45/19.09$. Our novelties and contributions thus lie in both the algorithm and the direction we investigate. The latter specifically leads to simpler solutions but higher accuracy.

**The importance of learning for the downstream tasks.** We found 176 classes of LVIS validation set in ILSVRC. That is, the corresponding ImageNet images used by MO-SAICOS are already seen by the pre-trained detector's backbone. Surprisingly, as shown in Table 6, MOSAICOS still

Table 7. **Instance segmentation on LVIS v0.5.** We use images from IN + G as Table 1. + [58]: include the balanced loss.

|  | AP | $AP_r$ | $AP_c$ | $AP_f$ |
|---|---|---|---|---|
| Mask R-CNN [24] | 24.38 | 15.98 | 23.96 | 28.27 |
| BaGS [47] | 26.25 | 17.97 | 26.91 | 28.74 |
| BALMS [58] | 27.00 | 19.60 | **28.90** | 27.50 |
| MOSAICOS | 26.26 | 19.63 | 26.60 | 28.49 |
| MOSAICOS + [58] | **27.86** | **20.44** | 28.82 | **29.62** |

Table 8. **LVIS v1.0 Results.** We report both object detection and instance segmentation performances of our method.

| OD Results | $AP^b$ | $AP^b_{50}$ | $AP^b_{75}$ | $AP^b_r$ | $AP^b_c$ | $AP^b_f$ |
|---|---|---|---|---|---|---|
| Faster R-CNN⋆ | 22.01 | 36.36 | 23.14 | 10.57 | 20.09 | 29.18 |
| MOSAICOS | **23.90** | **38.61** | **25.32** | **15.45** | **22.39** | **29.30** |
| IS Results | AP | $AP_{50}$ | $AP_{75}$ | $AP_r$ | $AP_c$ | $AP_f$ |
| Mask R-CNN | 22.59 | 35.44 | 23.87 | 12.31 | 21.30 | 28.55 |
| MOSAICOS | **24.49** | **38.02** | **25.87** | **18.30** | **23.00** | **28.87** |

leads to a notable gain for these classes, which not only justifies its efficacy, but also suggests the importance of learning the downstream tasks directly with those images.

### 6.4. Results on Instance Segmentation & LVIS v1.0

**Instance segmentation.** We also validate our approach on instance segmentation, in a similar manner: we prepare pseudo scene-centric images *with box labels* and use them in the first fine-tuning stage by optimizing the losses in Equation 1. That is, we do not optimize segmentation losses in this stage. We apply Mask R-CNN [29] with ResNet-50 as the backbone. Table 7 shows the results: *the baselines do not use extra object-centric images*. We see a notable gain against vanilla Mask R-CNN for rare and common classes, even if we have no segmentation labels on object-centric images. This supports the claim in [76]: even for detection and segmentation, the long-tailed problem is mainly in the classification sub-network. We perform on par with the state-of-the-art methods. Details are in the supplementary.

**LVIS v1.0 results.** We highlight consistent empirical results on LVIS v1.0, where our approach wins in both object detection and instance segmentation, using ResNet-50-FPN (see Table 8). More comparisons are in the supplementary.

## 7. Discussion and Conclusion

We investigate the use of object-centric images to facilitate long-tailed object detection on scene-centric images. We propose a concrete framework for this idea that is both simple and effective. Our results are encouraging, improving the baseline by a large margin on not only detecting but also segmenting rare objects. We hope that our study can attract more attention in using these already available but less explored object-centric images to overcome the long-tailed problem. Please see the supplementary for more discussion.

# References

[1] Flickr images. https://flickr.com/. 8

[2] Google images. https://www.google.com/imghp?hl=EN. 5, 8

[3] Hakan Bilen and Andrea Vedaldi. Weakly supervised deep detection networks. In *CVPR*, 2016. 2, 4

[4] Alexey Bochkovskiy, Chien-Yao Wang, and Hong-Yuan Mark Liao. Yolov4: Optimal speed and accuracy of object detection. *arXiv preprint arXiv:2004.10934*, 2020. 1, 2, 7, 12

[5] Matthew Brown and David G Lowe. Automatic panoramic image stitching using invariant features. *IJCV*, 74(1):59–73, 2007. 5

[6] Mateusz Buda, Atsuto Maki, and Maciej A Mazurowski. A systematic study of the class imbalance problem in convolutional neural networks. *Neural Networks*, 106:249–259, 2018. 2, 5

[7] Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Arechiga, and Tengyu Ma. Learning imbalanced datasets with label-distribution-aware margin loss. In *NeurIPS*, 2019.

[8] Nitesh V Chawla, Kevin W Bowyer, Lawrence O Hall, and W Philip Kegelmeyer. Smote: synthetic minority oversampling technique. *Journal of artificial intelligence research*, 16:321–357, 2002. 2

[9] Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jiarui Xu, Zheng Zhang, Dazhi Cheng, Chenchen Zhu, Tianheng Cheng, Qijie Zhao, Buyu Li, Xin Lu, Rui Zhu, Yue Wu, Jifeng Dai, Jingdong Wang, Jianping Shi, Wanli Ouyang, Chen Change Loy, and Dahua Lin. MMDetection: Open mmlab detection toolbox and benchmark. *arXiv preprint arXiv:1906.07155*, 2019. 18

[10] Yukang Chen, Peizhen Zhang, Zeming Li, Yanwei Li, Xiangyu Zhang, Gaofeng Meng, Shiming Xiang, Jian Sun, and Jiaya Jia. Stitcher: Feedback-driven data provider for object detection. *arXiv preprint arXiv:2004.12432*, 2020. 2, 7, 12

[11] Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. Class-balanced loss based on effective number of samples. In *CVPR*, 2019. 2

[12] Achal Dave, Piotr Dollár, Deva Ramanan, Alexander Kirillov, and Ross Girshick. Evaluating large-vocabulary object detectors: The devil is in the details. *arXiv preprint arXiv:2102.01066*, 2021. 15

[13] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 1, 2, 3, 4, 5, 8, 16

[14] Santosh K Divvala, Ali Farhadi, and Carlos Guestrin. Learning everything about anything: Webly-supervised visual concept learning. In *CVPR*, 2014. 2, 4

[15] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *IJCV*, 88(2):303–338, 2010. 2

[16] Li Fei-Fei, Rob Fergus, and Pietro Perona. Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories. In *CVPR*, 2004. 4

[17] Pedro F Felzenszwalb, Ross B Girshick, David McAllester, and Deva Ramanan. Object detection with discriminatively trained part-based models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(9):1627–1645, 2009. 1

[18] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *JMLR*, 17(1):2096–2030, 2016. 3, 16

[19] Jiyang Gao, Jiang Wang, Shengyang Dai, Li-Jia Li, and Ram Nevatia. Note-rcnn: Noise tolerant ensemble rcnn for semi-supervised object detection. In *ICCV*, 2019. 2, 4

[20] Ross Girshick. Fast R-CNN. In *ICCV*, 2015. 6

[21] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *CVPR*, 2014. 1

[22] Gregory Griffin, Alex Holub, and Pietro Perona. Caltech-256 object category dataset. 2007. 4

[23] Yandong Guo, Lei Zhang, Yuxiao Hu, Xiaodong He, and Jianfeng Gao. Ms-celeb-1m: A dataset and benchmark for large-scale face recognition. In *ECCV*, 2016. 2

[24] Agrim Gupta, Piotr Dollar, and Ross Girshick. LVIS: A dataset for large vocabulary instance segmentation. In *CVPR*, 2019. 1, 2, 3, 5, 6, 8, 16, 17, 18, 19

[25] James Hays and Alexei A Efros. Scene completion using millions of photographs. *ACM Transactions on Graphics (TOG)*, 26(3):4–es, 2007. 5

[26] Haibo He and Edwardo A Garcia. Learning from imbalanced data. *IEEE Transactions on knowledge and data engineering*, 21(9):1263–1284, 2009. 2

[27] Haibo He and Edwardo A Garcia. Learning from imbalanced data. *IEEE Transactions on knowledge and data engineering*, 21(9):1263–1284, 2009. 5

[28] Kaiming He, Ross Girshick, and Piotr Dollár. Rethinking imagenet pre-training. In *ICCV*, 2019. 2

[29] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask R-CNN. In *ICCV*, 2017. 1, 2, 6, 8, 16, 18

[30] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 1, 6, 13, 18

[31] Judy Hoffman, Sergio Guadarrama, Eric S Tzeng, Ronghang Hu, Jeff Donahue, Ross Girshick, Trevor Darrell, and Kate Saenko. Lsda: Large scale detection through adaptation. In *NIPS*, 2014. 2

[32] Judy Hoffman, Deepak Pathak, Trevor Darrell, and Kate Saenko. Detector discovery in the wild: Joint multiple instance and representation learning. In *CVPR*, 2015.

[33] Judy Hoffman, Deepak Pathak, Eric Tzeng, Jonathan Long, Sergio Guadarrama, Trevor Darrell, and Kate Saenko. Large scale visual recognition through adaptation using joint representation and multiple instance learning. *JMLR*, 17(1):4954–4984, 2016. 2

[34] Xinting Hu, Yi Jiang, Kaihua Tang, Jingyuan Chen, Chunyan Miao, and Hanwang Zhang. Learning to segment the tail. In *CVPR*, 2020. 2, 17

[35] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *CVPR*, 2017. 1

[36] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015. 4

[37] Jisoo Jeong, Seungeui Lee, Jeesoo Kim, and Nojun Kwak.

Consistency-based semi-supervised learning for object detection. In *NeurIPS*, 2019. 2

[38] Bingyi Kang, Saining Xie, Marcus Rohrbach, Zhicheng Yan, Albert Gordo, Jiashi Feng, and Yannis Kalantidis. Decoupling representation and classifier for long-tailed recognition. In *ICLR*, 2020. 2, 4, 5, 18

[39] Byungju Kim and Junmo Kim. Adjusting decision boundary for class imbalanced learning. *IEEE Access*, 8:81674–81685, 2020. 5

[40] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *IJCV*, 123(1):32–73, 2017. 2

[41] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012. 1

[42] Jason Kuen, Federico Perazzi, Zhe Lin, Jianming Zhang, and Yap-Peng Tan. Scaling object detection by transferring classification weights. In *ICCV*, 2019. 2

[43] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Malloci, Alexander Kolesnikov, et al. The open images dataset v4. *IJCV*, pages 1–26, 2020. 1

[44] Dong-Hyun Lee. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. In *ICML Workshop on challenges in representation learning*, 2013. 6

[45] Dong Li, Jia-Bin Huang, Yali Li, Shengjin Wang, and Ming-Hsuan Yang. Weakly supervised object localization with progressive domain adaptation. In *CVPR*, 2016. 2, 4

[46] Yandong Li, Di Huang, Danfeng Qin, Liqiang Wang, and Boqing Gong. Improving object detection with selective self-supervised self-training. In *ECCV*, 2020. 2, 4

[47] Yu Li, Tao Wang, Bingyi Kang, Sheng Tang, Chunfeng Wang, Jintao Li, and Jiashi Feng. Overcoming classifier imbalance for long-tail object detection with balanced group softmax. In *CVPR*, 2020. 1, 2, 7, 8, 17, 18, 19

[48] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *CVPR*, 2017. 1, 6

[49] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *ICCV*, 2017. 1

[50] 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, 3, 17

[51] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *ECCV*, 2016. 1

[52] 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 pretraining. In *ECCV*, 2018. 5

[53] George A Miller. WordNet: A lexical database for english. *Communications of the ACM*, 38(11):39–41, 1995. 3, 5, 18

[54] Constantine P Papageorgiou, Michael Oren, and Tomaso Poggio. A general framework for object detection. In *ICCV*, 1998. 1

[55] Vignesh Ramanathan, Rui Wang, and Dhruv Mahajan. DLWL: Improving detection for lowshot classes with weakly labelled data. In *CVPR*, 2020. 2, 4, 6, 7, 8, 12, 16

[56] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *CVPR*, 2016. 1

[57] Joseph Redmon and Ali Farhadi. Yolo9000: better, faster, stronger. In *CVPR*, 2017. 2

[58] Jiawei Ren, Cunjun Yu, Shunan Sheng, Xiao Ma, Haiyu Zhao, Shuai Yi, and Hongsheng Li. Balanced meta-softmax for long-tailed visual recognition. In *NeurIPS*, 2020. 2, 7, 8, 17

[59] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *NeurIPS*, 2015. 1, 3, 6, 13, 16

[60] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *IJCV*, 115(3):211–252, 2015. 2, 3, 5, 6, 18

[61] Shuai Shao, Zeming Li, Tianyuan Zhang, Chao Peng, Gang Yu, Xiangyu Zhang, Jing Li, and Jian Sun. Objects365: A large-scale, high-quality dataset for object detection. In *ICCV*, 2019. 1, 3

[62] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. 1

[63] Kihyuk Sohn, Zizhao Zhang, Chun-Liang Li, Han Zhang, Chen-Yu Lee, and Tomas Pfister. A simple semi-supervised learning framework for object detection. *arXiv preprint arXiv:2005.04757*, 2020. 2

[64] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *CVPR*, 2015. 1

[65] Jingru Tan, Xin Lu, Gang Zhang, Changqing Yin, and Quanquan Li. Equalization loss v2: A new gradient balance approach for long-tailed object detection. In *CVPR*, 2021. 16, 17, 18

[66] Jingru Tan, Changbao Wang, Buyu Li, Quanquan Li, Wanli Ouyang, Changqing Yin, and Junjie Yan. Equalization loss for long-tailed object recognition. In *CVPR*, 2020. 1, 2, 4, 16, 17, 18

[67] Jingru Tan, Gang Zhang, Hanming Deng, Changbao Wang, Lewei Lu, Quanquan Li, and Jifeng Dai. 1st place solution of lvis challenge 2020: A good box is not a guarantee of a good mask. *arXiv preprint arXiv:2009.01559*, 2020. 17

[68] Kaihua Tang, Jianqiang Huang, and Hanwang Zhang. Long-tailed classification by keeping the good and removing the bad momentum causal effect. In *NeurIPS*, 2020. 2

[69] Yuxing Tang, Josiah Wang, Boyang Gao, Emmanuel Dellandréa, Robert Gaizauskas, and Liming Chen. Large scale semi-supervised object detection using visual and semantic knowledge transfer. In *CVPR*, 2016. 2

[70] Yuxing Tang, Josiah Wang, Xiaofang Wang, Boyang Gao, Emmanuel Dellandréa, Robert Gaizauskas, and Liming Chen. Visual and semantic knowledge transfer for large scale semi-supervised object detection. *TPAMI*, 40(12):3045–

3058, 2017. 2

[71] Bart Thomee, David A Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. YFCC100M: The new data in multimedia research. *Communications of the ACM*, 59(2):64–73, 2016. 1, 2, 6, 7, 12

[72] Jasper Uijlings, Stefan Popov, and Vittorio Ferrari. Revisiting knowledge transfer for training object class detectors. In *CVPR*, 2018. 2, 4

[73] 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. 1, 2

[74] Paul Viola, Michael Jones, et al. Robust real-time object detection. *IJCV*, 4(34-47):4, 2001. 1

[75] Jiaqi Wang, Wenwei Zhang, Yuhang Zang, Yuhang Cao, Jiangmiao Pang, Tao Gong, Kai Chen, Ziwei Liu, Chen Change Loy, and Dahua Lin. Seesaw loss for long-tailed instance segmentation. In *CVPR*, 2021. 16, 17, 18, 19

[76] Tao Wang, Yu Li, Bingyi Kang, Junnan Li, Junhao Liew, Sheng Tang, Steven Hoi, and Jiashi Feng. The devil is in classification: A simple framework for long-tail instance segmentation. In *ECCV*, 2020. 2, 8, 17

[77] Xin Wang, Thomas E Huang, Trevor Darrell, Joseph E Gonzalez, and Fisher Yu. Frustratingly simple few-shot object detection. In *ICML*, 2020. 1, 2, 7, 17

[78] Zirui Wang, Zihang Dai, Barnabás Póczos, and Jaime Carbonell. Characterizing and avoiding negative transfer. In *CVPR*, 2019. 2, 3

[79] Chen Wei, Kihyuk Sohn, Clayton Mellina, Alan Yuille, and Fan Yang. Crest: A class-rebalancing self-training framework for imbalanced semi-supervised learning. *arXiv preprint arXiv:2102.09559*, 2021. 15

[80] Jialian Wu, Liangchen Song, Tiancai Wang, Qian Zhang, and Junsong Yuan. Forest R-CNN: Large-vocabulary long-tailed object detection and instance segmentation. In *ACM MM*, 2020. 2, 17

[81] Yuxin Wu, Alexander Kirillov, Francisco Massa, Wan-Yen Lo, and Ross Girshick. Detectron2. https://github.com/facebookresearch/detectron2, 2019. 6, 18

[82] Qizhe Xie, Minh-Thang Luong, Eduard Hovy, and Quoc V Le. Self-training with noisy student improves imagenet classification. In *CVPR*, 2020. 15

[83] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *CVPR*, 2017. 18

[84] Han-Jia Ye, Hong-You Chen, De-Chuan Zhan, and Wei-Lun Chao. Identifying and compensating for feature deviation in imbalanced deep learning. *arXiv preprint arXiv:2001.01385*, 2020. 4, 5

[85] Fisher Yu, Haofeng Chen, Xin Wang, Wenqi Xian, Yingying Chen, Fangchen Liu, Vashisht Madhavan, and Trevor Darrell. BDD100K: A diverse driving dataset for heterogeneous multitask learning. In *CVPR*, 2020. 1

[86] Amir R Zamir, Alexander Sax, William Shen, Leonidas J Guibas, Jitendra Malik, and Silvio Savarese. Taskonomy: Disentangling task transfer learning. In *CVPR*, 2018. 2, 3

[87] Fan Zhang and Feng Liu. Parallax-tolerant image stitching. In *CVPR*, 2014. 5

[88] Yubo Zhang, Pavel Tokmakov, Martial Hebert, and Cordelia Schmid. A study on action detection in the wild. *arXiv preprint arXiv:1904.12993*, 2019. 2

[89] Dongzhan Zhou, Xinchi Zhou, Hongwen Zhang, Shuai Yi, and Wanli Ouyang. Cheaper pre-training lunch: An efficient paradigm for object detection. In *ECCV*. Springer. 2, 7

[90] Xiangxin Zhu, Dragomir Anguelov, and Deva Ramanan. Capturing long-tail distributions of object subcategories. In *CVPR*, 2014. 1

[91] Barret Zoph, Golnaz Ghiasi, Tsung-Yi Lin, Yin Cui, Hanxiao Liu, Ekin Dogus Cubuk, and Quoc Le. Rethinking pre-training and self-training. In *NeurIPS*, 2020. 1, 2, 6, 7, 15

# Supplementary Material

In this Supplementary Material, we provide details and results omitted in the main text.

## A. Contribution and Novelty

Our main contributions are in the idea of using object-centric images (OCI) to facilitate long-tailed object detection on scene-centric images (SCI) as well as a concrete implementation of this idea that is both simple and effective. This is by no means trivial; for instance, a related work [55] with a more sophisticated approach can hardly improve the accuracy (Table 5 of the main paper). While most existing works focus on *designing new algorithms* to learn from long-tailed data, our proposal is orthogonal to them, and can be combined together for further improvement.

Although leveraging auxiliary data to improve *common* object detection has been studied previously, existing works typically assume access to well prepared data from a similar visual domain, with sufficient object instances. However, collecting and annotating such auxiliary data is extremely challenging in *long-tailed* object detection. In contrast, our method does not have such a limitation as we make use of *object-centric* images readily available over the Internet (via search engines), which contains sufficient object instances though in a slight different domain. Particularly, we observe that making use of such rich *object-centric* images (from



Figure 6. **Different stitching methods.** MOSAICOS introduces *more diverse* examples by leveraging object-centric images, while existing methods [4, 10] only perform data augmentation using scene-centric images.

ImageNet) leads to more superior empirical performances against [55], which uses YFCC-100M [71].

To enable more general applicability, we make the design of our framework as straightforward as possible. Along this process, two challenges are identified, *i.e.*, the gap between visual domains and the lack of object labels. To address them, we investigate simple algorithms such as fixed box locations, mosaicking, and multi-stage training. We note that more sophisticated techniques can be incorporated as well. The facts that (a) *our framework performs on par with state-of-the-art long-tailed detection methods* and (b) *many existing techniques can be easily plugged into our framework* further justify the potential of this promising direction.

While several components of our framework — mosaic, pseudo-labeling, two-stage fine-tuning — have been individually explored in prior works in different contexts, *a suitable combination is essential and novel* for our idea to work. Further, our use of mosaic on OCI is different from [4, 10], as shown in Figure 6. Our contributions also include extensive analysis that justifies the importance of each component. These insights led to a simple and effective framework, which we consider a strength. For example, our LORE approach (§ B.2) could have provided methodological novelty. But its small gain over simple fixed locations does not justify the inclusion of it into our final framework.

## B. Pseudo-Label Generation

### B.1. Trust the calibrated detector and image labels

We provide analysis on pseudo-label generation with detector calibration and imputation using image class labels.

**Detector calibration.** As mentioned in § 5.1 of the main paper, we calibrate the pre-trained detector by assigning each class a different confidence threshold according to the class size — rare classes have lower thresholds. Figure 7 illustrates the difference with and without detector calibration, and with and without imputation using the image class labels. By assigning each class a different confidence threshold, the calibrated detector outputs more detected boxes, indicating that many rare and common objects are missed by the pre-trained detector due to low confidence scores (Figure 7 (a) vs. (c)). However, simply applying calibration can hardly correct the wrong labels that have already been biased toward the frequent classes (blue boxes in Fig-
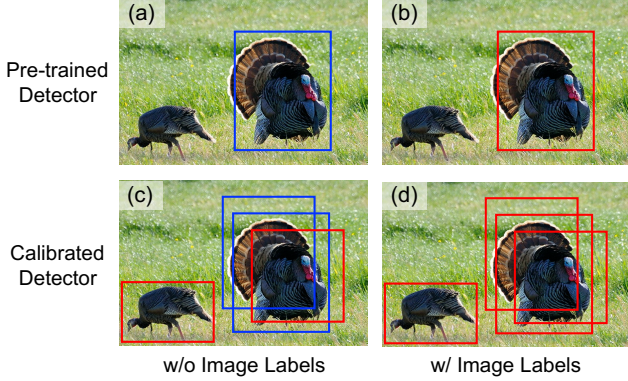
Figure 7. **A comparison of pseudo-label generation with detector calibration and imputation using image class labels.** (a) trust the detector (D), (b) trust the detector + image class labels (D†), (c) trust the calibrated detector, and (d) trust the calibrated detector + class image labels (D‡). The image label is "turkey", a rare class in LVIS. Red/Blue boxes are labeled as "turkey"/other classes. See § 5.1 of the main paper for details.

ure 7 (c)). Next, we explore the idea of bringing the best of image class labels to correct noisy detected labels.

**The importance of imputation with image class labels.** For object-centric images, most of the object instances belong to the image's class label. We therefore improve "trust the pre-trained detector" (Figure 7 (a)) and "trust the calibrated detector" (Figure 7 (c)) by assigning each box the image class label (see Figure 7 (b) and (d)). As shown in Figure 7 and Table 9, we see significant improvements for both the pre-trained and calibrated detectors. Specifically, assigning the image class label for each box can largely boost the performance for rare objects ($AP^b_r$).

## B.2. Details on LORE

Figure 8 shows the pipeline of localization by region removal (LORE), which is introduced in § 5.1 of the main paper for pseudo-label generation. Concretely, LORE takes an object-centric image as the input and identifies the locations of the target object (*i.e.*, that of image label) in the image. The whole pipeline consists of three major components: (1) classifier training, (2) box pre-filtering, and (3) localization by removal. We describe each step as follows.

**Classifier training.** We train a ResNet-50 [30] image classifier with all object-centric images. For LVIS v0.5 dataset, we follow the conventional training procedure[8] to train a $1,230$-way ResNet classifier. Specifically, we train the networks with 90 epoch and achieve $74\%$ top-1 training accuracy. We use this pre-trained classifier to rank object regions in object-centric images.

**Box pre-filtering.** We feed an object-centric image into

Table 9. **Results with different pseudo-labels.** We use ImageNet-21K as the source of object-centric images and report the results of object detection on LVIS v0.5 val. **Detector**: object detector used for generating pseudo-label bounding boxes; **CL:** assign each box the image Class Label instead of the predicted class label.

| | Detector | CL | $AP^b$ | $AP^b_r$ | $AP^b_c$ | $AP^b_f$ |
|---|---|---|---|---|---|---|
| Faster R-CNN⋆ | – | – | 23.35 | 12.98 | 22.60 | 28.42 |
| MOSAICOS | Pre-trained | ✗ | 23.04 | 13.93 | 21.51 | 28.14 |
| | Pre-trained | ✓ | 24.66 | 17.45 | **23.62** | 28.83 |
| | Calibrated | ✗ | 24.03 | 13.13 | 23.51 | **29.04** |
| | Calibrated | ✓ | **24.93** | **19.31** | 23.51 | 28.95 |

the pre-trained *object detector* and collect detection results. Concretely, we take the top 300 detected boxes of Faster R-CNN [59] and drop each box's predicted class label. Next, we apply non-maximum suppression (NMS) over all the 300 boxes using a threshold of 0.5 to remove highly-overlapped ones. Basically, we trust the detected box locations (*i.e.*, they do contain objects), but will recheck which of them belongs to the target object.

To further reduce the number of candidate boxes, we sort the boxes by their initial detection confidence (in the descending order) and then remove the corresponding regions from the image *in turn*[9], every time followed by applying the image classifier to the resulting image. We stop this process until the classification confidence of the target class goes below a certain threshold. We then collect the removed box locations, which together have likely covered the target objects (high recall, but likely low precision), to be the candidate box pool for the next step.

**Localization by removal.** To accurately identify which candidate truly belongs to the target class, we *re-rank* the candidates by how much removing each boxed region *alone* reduces the image classifier's confidence on the target class. We then follow the descending order to remove these boxed regions *in turn* until the classifier fail to predict the target class or the *confidence reducing ratio*[10] achieves a certain threshold. Finally, the bounding boxes of the removed regions are collected as the pseudo ground-truths for the image. More examples can be found in Figure 9.

## B.3. Discussion on fixed locations vs. LORE

Both fixed locations and LORE use accurate image class labels. Even though LORE gives more accurate object locations (see Figure 9) in pseudo-label generation, the resulting detector with fixed location is just slightly worse than that with LORE. We attribute this small gap partially to two-stage fine-tuning, which adapts the detector back to accurately labeled scene-centric images. As shown in Table 10,

---

[8]https://github.com/pytorch/examples/tree/master/imagenet

[9]We crop out the corresponding image regions and replacing them by gray-color patches.

[10]We define the *confidence reducing ratio* as the relative confidence drop on the target class label before and after removing boxes.

13

Figure 8. **Illustration of LORE.** We first apply a pre-trained detector to obtain candidate boxes, followed by pre-filtering. We then sort the remaining boxes using an image classifier. Finally, we remove the boxes in turn until the classifier fail to predict the target image label. The numbers at image corners indicate the confidence reducing ratio. Negative values mean the confidence increases after removing outliers.



Figure 9. **Box locations of different pseudo-label generation methods.** We show (a) fixed locations, (b) trust the pre-trained detector, (c) trust the calibrated detector, and (d) localization by region removals (LORE). The green boxes are the pseudo ground-truth locations found on each object-centric image alone before multiple images are stitched together. We can see that LORE accurately locates the target object in each sub-image while detection results are much noisy. Image class labels are listed on the corner of each sub-image in column (a).

14

Table 10. **Fixed locations vs. LORE.** We report object detection results on LVIS v0.5 val. **P-GT**: ways to generate pseudo-labels.

|  | P-GT | $AP^b$ | $AP^b_r$ | $AP^b_c$ | $AP^b_f$ |
|---|---|---|---|---|---|
| Single-stage | Fixed | 20.09 | 12.96 | 19.08 | 24.20 |
|  | LORE | 21.44 | 14.95 | 20.74 | 24.91 |
| MOSAICOS | Fixed | 24.75 | 19.73 | 23.44 | 28.39 |
|  | LORE | 24.83 | 20.06 | 23.25 | 28.71 |

LORE notably surpasses fixed locations if we apply single-stage fine-tuning.

## B.4. Discussion on pseudo-label generation

In this subsection, we discuss multiple ways for generating pseudo-labels in object-centric images. From the viewpoint of teacher models (*i.e.,* the pre-trained detector learned from a long-tailed distribution), we found that (1) the pre-trained detector is biased toward head classes, missing many accurate rare class predictions which have lower confidence scores; (2) detector calibration is useful to discover more bounding boxes for rare and common objects but can hardly correct wrong predicted labels. Our observations share the similar insights with a recent study [12] on large-vocabulary object detection.

From the other viewpoint of fine-tuning with pseudo scene-centric images, we found that imputation using image class labels leads to a notable performance gain regardless of inaccurate box locations (*e.g.*, fixed box locations). This is probably due to two reasons. First, dense boxes (like six fixed locations) can be treated as data augmentation for training the object detector. Second, our two-stage fine-tuning is beneficial in learning with noisy data, *i.e.,* first on noisy pseudo scene-centric images and then on the clean labeled data from LVIS.

Other possibilities for pseudo-label generation include (1) iteratively improving the teacher detector by noisy student learning [82] and (2) calibrating the detector with more advanced approaches for class-imbalanced semi-supervised learning [79], etc.

## C. Additional Ablation on Image Mosaicking

*Does mosaicking more images help?* In this section, we investigate the effect of different types of layouts for stitching object-centric images, *i.e.*, $1 \times 1$ (which is the original object-centric image), $2 \times 2$ mosaic, and $3 \times 3$ mosaic. We evaluate them under the same experimental settings: we use ImageNet-21K as the source of object-centric images ($1,016$ classes) and stitch images from the same class and use the 6 fixed locations as pseudo ground-truths. Table 11 shows the comparison of object detection results on LVIS v0.5 dataset. We see that $2 \times 2$ and $3 \times 3$ mosaics perform similarly and both outperform the $1 \times 1$ OCI (on $AP^b$ and $AP^b_r$). An example with different layouts of $2 \times 2$ and $3 \times 3$ mosaics is shown in Figure 10.



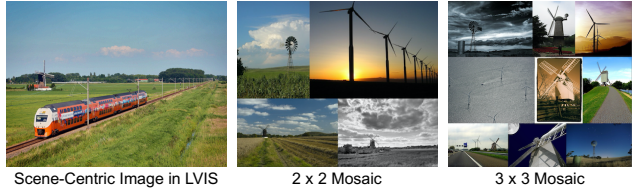Scene-Centric Image in LVIS    2 x 2 Mosaic    3 x 3 Mosaic

Figure 10. **Different layouts of mosaics.** We show different types of mosaics from the same category ("windmill"). The $2 \times 2$ mosaic image (middle) and the real scene-centric image (left) in the LVIS dataset look alike in terms of appearance and structure while the $3 \times 3$ mosaic image (right) is much crowded.

Table 11. **Comparison of different types of mosaic images.** Here we use ImageNet-21K as the source of object-centric images and stitch images from the *same* class and use the 6 fixed locations as pseudo ground-truths. $1 \times 1$ OCI means directly using the original object-centric images. Results are reported on LVIS v0.5 val. We can see that $2 \times 2$ mosaic gives better performance on all classes. The best result per column is in bold font.

|  | $AP^b$ | $AP^b_r$ | $AP^b_c$ | $AP^b_f$ |
|---|---|---|---|---|
| Faster R-CNN⋆ | 23.35 | 12.98 | 22.60 | 28.42 |
| $1 \times 1$ OCI | 24.27 | 16.97 | **23.29** | **28.42** |
| $3 \times 3$ Mosaic | 24.29 | 18.14 | 23.13 | 28.21 |
| $2 \times 2$ Mosaic | **24.48** | **18.76** | 23.26 | 28.29 |

## D. Further Analysis on Self-training

We show detailed comparison results of self-training baseline in Table 12 to further demonstrate the effectiveness of the mosaicking and two-stage fine-tuning in our MO-SAICOS framework. We follow the self-training method with the normalization loss in [91].

**Mosaicking is also beneficial for self-training.** We first study the vanilla self-training that directly learns object-centric (without mosaicking) and scene-centric images jointly. Specifically, we apply the pre-trained detector to generate pseudo-labels on the object-centric images (D). Next, the pre-trained detector is trained to jointly optimize the losses on human labels from LVIS and pseudo labels on object-centric images. We compare with and without image mosaic in Table 12: image mosaicking improves $AP^b/AP^b_r$ from $22.00/14.04$ to $22.71/14.52$, demonstrating the effectiveness of mosaicking object-centric images to mitigate the domain discrepancy between two types of images.

**Self-training vs. our two-stage fine-tuning.** To further improve the performance of self-training, we apply "trusted the calibrated detector + image class labels" (D‡) as the pseudo-labeling method, which leads to a much higher detection accuracy than "trusted the pre-trained detector" (D) for our MOSAICOS (cf. Table 1 of the main paper and the last row vs. the first row in Table 12). With this pseudo-labeling method, we see a notable gain against "trust the pre-trained detector" (D) for self-training.

We further compare the self-training procedure that fine-

Table 12. **Comparison to self-training. Mosaic:** ✓means 2×2 image mosaicking from different classes. **P-GT:** ways to generate pseudo-labels (**D**: trust the pre-trained detector, **D‡**: trust the calibrated detector and image class label).

| | Mosaic | P-GT | $AP^b$ | $AP^b_r$ | $AP^b_c$ | $AP^b_f$ |
|---|---|---|---|---|---|---|
| Faster R-CNN⋆ | – | – | 23.35 | 12.98 | 22.60 | 28.42 |
| Self-training | ✗ | D | 22.00 | 14.04 | 20.41 | 27.18 |
| | ✓ | D | 22.71 | 14.52 | 21.41 | 27.61 |
| | ✓ | D‡ | 23.65 | 16.30 | 22.55 | 27.96 |
| MOSAICOS | ✓ | D | 23.04 | 13.93 | 21.51 | 28.14 |
| | ✓ | D‡ | 24.93 | 19.31 | 23.51 | 28.95 |

Table 13. **Comparison to adversarial training.** Results are reported on LVIS v0.5 validation set.

| | Mosaic | $AP^b$ | $AP^b_r$ | $AP^b_c$ | $AP^b_f$ |
|---|---|---|---|---|---|
| Single-stage | ✓ | 20.09 | 12.96 | 19.08 | 24.20 |
| Adv. training | ✗ | 20.92 | 11.42 | 19.23 | 26.85 |
| | ✓ | 22.87 | 14.48 | 22.02 | 27.28 |
| MOSAICOS | ✗ | 24.27 | 16.97 | 23.29 | **28.42** |
| | ✓ | **24.75** | **19.73** | **23.44** | 28.39 |

tunes the detector simultaneously with object-centric and scene-centric images to our MOSAICOS with two-stage fine-tuning (again in Table 12). MOSAICOS outperforms self-training (with either D or D‡) in most metrics, demonstrating the strength of two-stage fine-tuning which first learns with object-centric images and then scene-centric images. This two-stage pipeline is not only robust to noisy pseudo scene-centric data but also able to tie the detector to its final application domain with real scene-centric images.

## E. Data Quality of Object-Centric Images

Our main results are based on the ImageNet dataset [13]. We included Google/Flickr images (Table 5 in the main text) mainly to analyze the effect of data quality and compare to [55]. As shown in Figure 11, most Google images searched by object names are object-centric, even for those not ranked on the top. Following the experimental setup in Table 5 of the main paper, we further experiment with 500 Google images per class: $AP^b$ is improved from 24.45 to 24.63. For images that are less object-centric, LORE can give better pseudo-labels than the fixed heuristic; our two-stage fine-tuning is robust to noise. Moreover, there are extensive works on de-noising web data that we can leverage to further improve our scalability and applicability. That being said, we neither focus on web images/crowd-sourcing nor suggest that human efforts (*e.g.*, ImageNet) are not needed. Our claim is that rare objects that are hard to collect from SCI are easier to collect from OCI, which opens up a new way to tackle long-tailed object detection.

## F. Comparison to Adversarial Training

We apply adversarial training (*e.g.*, [18]) to jointly train the detector with LVIS and pseudo scene-centric images. Concretely, we train an additional domain classifier to differentiate the LVIS images and pseudo scene-centric images, and incorporate a gradient reversal layer (GRL) [18] to minimize the discrepancy between their features to overcome the domain gap. We show comparisons in Table 13. Adversarial training outperforms naive joint (*i.e.*, single-stage) training, and MOSAICOS (with two-stage fine-tuning using each source) surpasses adversarial training.

## G. Implementation Details of MOSAICOS

### G.1. Details on object detection

As mentioned in § 6.1 of the main paper, we use Faster R-CNN [59] as our base detector and further extend the training process with another 90K iterations and select the checkpoint with the best $AP^b$ as Faster R-CNN⋆. We use Faster R-CNN⋆ as our main baseline to ensure that the improvement of MOSAICOS does not simply come from training (*i.e.*, fine-tuning) with more epochs.

For MOSAICOS, we first fine-tune Faster R-CNN⋆ with pseudo scene-centric images, and then fine-tune it with the LVIS training set again. Both stages are trained end-to-end with stochastic gradient descent with all training losses in Equation 1 of the main paper, using a mini-batch size of 16, momentum of 0.9, weight decay of $10^{-4}$, and learning rate of $2 \times 10^{-4}$. *Unlike other long-tailed methods [65, 66, 75][11], there is no additional hyper-parameter in our framework.*

### G.2. Details on instance segmentation

**Background on instance segmentation.** We apply Mask R-CNN [29], which adopts the two-stage network architecture similar to Faster R-CNN [59], with an identical first stage RPN. In the second stage, in addition to predicting the class label and box offset, Mask R-CNN further outputs a binary segmentation mask for each proposal. Formally, during training, the entire Mask R-CNN is learned with four loss terms

$$\mathcal{L} = \mathcal{L}_{rpn} + \mathcal{L}_{cls} + \mathcal{L}_{reg} + \mathcal{L}_{mask}, \qquad (2)$$

where the RPN loss $\mathcal{L}_{rpn}$, classification loss $\mathcal{L}_{cls}$, and box regression loss $\mathcal{L}_{reg}$ are identical to those defined in [59]. The mask loss $\mathcal{L}_{mask}$ is learned via an average binary cross-entropy objective.

**Multi-stage training for instance segmentation.** We first train a Mask R-CNN using labeled scene-centric images from LVIS with instance segmentation annotations [24]. All the fours loss terms in Equation 2 are optimized.

We then fine-tune the model using the pseudo scene-centric images that are generated from object-centric images. We use these images (only with box pseudo-labels)

---
[11]Both EQL(v2) [65, 66] and Seesaw loss [75] introduce (multiple) additional hyper-parameters.

Figure 11. **Google Images for the most rare classes in LVIS.** We show the top 5 retrieved images and images ranked around 500.

Table 14. **Object detection on LVIS v0.5**. We use ImageNet + Google Images. MSCOCO: for pre-training. [58]: balanced loss. Within each column, red/blue indicates the best/second best.

|  | MSCOCO | [58] | $AP^b$ | $AP_r^b$ | $AP_c^b$ | $AP_f^b$ |
|---|---|---|---|---|---|---|
| RFS [24] |  |  | 23.35 | 12.98 | 22.60 | 28.42 |
| EQL [66] |  |  | 23.30 | – | – | – |
| LST [34] |  |  | 22.60 | – | – | – |
| BaGS [47] | ✓ |  | 25.96 | 17.65 | 25.75 | 29.54 |
| TFA [77] |  |  | 24.40 | 16.90 | 24.30 | 27.70 |
| MOSAICOS |  |  | 25.01 | 20.25 | 23.89 | 28.32 |
|  | ✓ |  | 26.28 | 17.37 | 26.13 | 30.02 |
|  |  | ✓ | 26.83 | 21.00 | 26.31 | 29.81 |
|  | ✓ | ✓ | 28.06 | 19.11 | 28.23 | 31.41 |

Table 15. **Instance segmentation on LVIS v0.5.** Our MOSAICOS uses images from ImageNet and Google Images. + [58]: include the balanced loss in the second stage fine-tuning. Within each column, red/blue indicates the best/second best.

|  | AP | $AP_r$ | $AP_c$ | $AP_f$ |
|---|---|---|---|---|
| RFS [24] | 24.38 | 15.98 | 23.96 | 28.27 |
| EQL [66] | 22.80 | 11.30 | 24.70 | 25.10 |
| LST [34] | 23.00 | – | – | – |
| SimCal [76] | 23.40 | 16.40 | 22.50 | 27.20 |
| Forest RCNN [80] | 25.60 | 18.30 | 26.40 | 27.60 |
| BaGS [47] | 26.25 | 17.97 | 26.91 | 28.74 |
| BALMS [58] | 27.00 | 19.60 | 28.90 | 27.50 |
| EQL v2 [65] | 27.10 | 18.60 | 27.60 | 29.90 |
| MOSAICOS | 26.26 | 19.63 | 26.60 | 28.49 |
| MOSAICOS + [58] | 27.86 | 20.44 | 28.82 | 29.62 |

to fine-tune the model using $\mathcal{L}_{cls}$, $\mathcal{L}_{rpn}$, and $\mathcal{L}_{reg}$. In other words, we do not optimize $\mathcal{L}_{mask}$. Any network parameters that affect $\mathcal{L}_{cls}$, $\mathcal{L}_{rpn}$, and $\mathcal{L}_{reg}$, especially those in the backbone feature network (except the batch-norm layers), can be updated.

After this stage, we fine-tune the whole network again with labeled scene-centric images from LVIS, using all the four loss terms in Equation 2. The training procedure and other implementation details for instance segmentation are exactly the same as object detection in § G.1.

## H. Experimental Results on LVIS v0.5

Due to space limitations, we only compared with state-of-the-art methods in Table 3 and Table 7 of the main paper. In this section, we provide detailed comparisons with more previous works on LVIS v0.5.

**Object detection on LVIS v0.5.** There are not many papers reporting detection results on LVIS. In Table 14, we further include EQL [66] and LST [34], together with BaGS [47] and TFA [77], as the compared methods. MOSAICOS outperforms all baselines except BaGS [47]. We note that, BaGS is pre-trained on COCO [50] while MOSAICOS is initialized from ResNet-50 that is pre-trained on ImageNet-1K (ILSVRC). By using the COCO pre-trained backbone as the initialization, MOSAICOS outperforms BaGS on nearly all metrics. Moreover, when combined with [58], MOSAICOS can further boost the state-of-the-art performance.

**Instance segmentation on LVIS v0.5.** The comparison re-

sults on LVIS 0.5 instance segmentation are presented in Table 15, including the baseline models with RFS [24] for re-sampling, EQL(v2) [65, 66] for re-weighting, LST [34] for incremental learning, SimCal [76] and BaGS [47] for de-coupled training, Forest R-CNN [80] for hierarchy classification, and BALMS [58] for a balanced softmax loss. MOSAICOS can perform on a par with or even better than the compared methods without any additional hyper-parameter tuning like in [65, 66, 75]. By combined with [58], MOSAICOS achieves the stat-of-the-art performance of $27.86/20.44$ AP/AP$_r$, showing the compatibility of MOSAICOS. We expect that MOSAICOS could be further improved by incorporating other long-tailed learning strategies [47, 58, 67, 75, 77].

## I. Experimental Results on LVIS v1.0

### I.1. Setup

**Dataset statistics.** We further evaluate MOSAICOS on LVIS v1.0 [24]. The total dataset size has been expanded to ∼160K images and ∼2M instance annotations. The total number of categories has decreased slightly (from 1,230 to 1,203) due to a more stringent quality control. More specifically, LVIS v1.0 adds 52 new classes while drops 79 classes from LVIS v0.5. The validation set has been expanded from 5K images to 20K images. Table 17 gives a summary of the statistics of the two versions of LVIS dataset. We follow

Table 16. **Number of overlapped classes in LVIS and ImageNet.** In LVIS and ImageNet, each category can be identifed by a unique WordNet synset ID. We match LVIS classes to ImageNet ones and show the number of the overlapped classes. Specifically, we show # LVIS classes / # overlapped to ImageNet-21K / # overlapped to ImageNet-1K (ILSVRC).

| Version | Split | Frequent | Common | Rare | Overall |
|---|---|---|---|---|---|
| v0.5 | Train | 315 / 253 / 85 | 461 / 387 / 96 | 454 / 385 / 71 | 1230 / 1025 / 252 |
| | Val | 313 / 252 / 21 | 392 / 329 / 84 | 125 / 106 / 71 | 830 / 678 / 176 |
| v1.0 | Train | 405 / 331 / 87 | 461 / 390 / 96 | 337 / 277 / 64 | 1203 / 998 / 247 |
| | Val | 405 / 331 / 87 | 452 / 382 / 92 | 178 / 144 / 37 | 1035 / 857 / 216 |

Table 17. **Statistics of LVIS v0.5 and v1.0 datasets.**

| Version | Type | Train | Val | Test |
|---|---|---|---|---|
| v0.5 | # Image | 57,263 | 5,000 | 19,761 |
| | # Class | 1,230 | 830 | - |
| | # Instance | 693,958 | 50,763 | - |
| v1.0 | # Image | 100,170 | 19,809 | 19,822 |
| | # Class | 1,203 | 1,035 | - |
| | # Instance | 1,270,141 | 244,707 | - |

Table 18. **Instance segmentation on LVIS v1.0.** We list multiple Mask R-CNN baselines whose accuracy are notably different due to differences in implementation, which may affect the accuracy of the corresponding proposed methods. MOSAICOS outperforms Mask R-CNN and many other methods on most of the metrics.

| Backbone | Method | AP | $AP_r$ | $AP_c$ | $AP_f$ |
|---|---|---|---|---|---|
| R-50 | Mask RCNN [24][†1] | 22.59 | 12.31 | 21.30 | 28.55 |
| | Mask RCNN [24][⋆2] | 23.70 | 13.50 | 22.80 | 29.30 |
| | Mask RCNN [24][§1] | 22.20 | 11.50 | 21.20 | 28.00 |
| | cRT [38][§1] | 22.10 | 11.90 | 20.20 | 29.00 |
| | BaGS [47][§1] | 23.10 | 13.10 | 22.50 | 28.20 |
| | EQL v2 [65][§1] | 23.70 | 14.90 | 22.80 | 28.60 |
| | EQL v2 [65][§2] | 25.50 | 17.70 | 24.30 | 30.20 |
| | Seesaw [75][⋆2] | 26.40 | 19.60 | 26.10 | 29.80 |
| | MOSAICOS [†1] | 24.49 | 18.30 | 23.00 | 28.87 |
| R-101 | Mask RCNN [24][†1] | 24.82 | 15.18 | 23.71 | 30.31 |
| | Mask RCNN [24][⋆2] | 25.50 | 16.60 | 24.50 | 30.60 |
| | EQL [66][⋆2] | 26.20 | 17.00 | 26.20 | 30.20 |
| | BaGS [47][⋆2] | 25.80 | 16.50 | 25.70 | 30.10 |
| | Seesaw [75][⋆2] | 28.10 | 20.00 | 28.00 | 31.90 |
| | MOSAICOS[†1] | 26.77 | 20.79 | 25.76 | 30.53 |
| X-101 | Mask RCNN [24][†1] | 26.62 | 17.51 | 25.51 | 31.86 |
| | MOSAICOS[†1] | 28.31 | 21.74 | 27.25 | 32.36 |

[†]: Our implementations with RFS [24].
[⋆]: Results reported in [75]. All models trained with RFS [24].
[§]: Results reported in [65].
[1]: 1x schedule. [2]: 2x schedule.

the experimental setups of LVIS v0.5 to use category synset ID [53] to search for the corresponding classes in ImageNet-21K dataset [60]. In total, we collect 753, 700 object-centric images. Table 16 shows the detailed statistics of the number of overlapped classes in those datasets. We also search 100 images for each class via Google Images.

**Our settings.** For instance segmentation, we use Mask R-CNN [29] with instance segmentation annotations. The training scheme is the same as that for Faster R-CNN in object detection. Specifically, we follow the default training configurations in [81] with 1x schedule[12].

For the MOSAICOS training (cf. § G.2), we first fine-tune the baseline Mask R-CNN for 90K iterations with pseudo scene-centric images using only box annotations. Our pseudo scene-centric images are synthesized with $2 \times 2$ mosaic from random classes of ImageNet-21K and Google images. We use the boxes with 6 fixed locations as pseudo ground-truths. After that, We end-to-end fine-tune the entire model for another 90K iterations using the LVIS training set with all four losses. The network parameters of the mask head are initialized by the baseline Mask RCNN model. Both two fine-tuning steps are trained with stochastic gradient descent with a mini-batch size of 16, momentum of 0.9, weight decay of $10^{-4}$, and learning rate of $2 \times 10^{-4}$.

### I.2. Instance segmentation on LVIS v1.0

Table 18 shows detailed results on instance segmentation. We mainly compare with Mask R-CNN and two recent papers [65, 75], which reported instance segmentation results and re-implemented some other methods on LVIS v1.0. We evaluate MOSAICOS with three different backbone models: ResNet-50 [30], ResNet-101 [30], and ResNeXt-101 [83]: MOSAICOS consistently outperforms the Mask R-CNN baseline especially for rare classes.

We note that, EQL v2 [65] and Seesaw loss [75] were implemented by a different framework [9] and reported results with a stronger 2x training schedule. *Thus, the accuracy gap between different methods may be partially affected by these factors.* This can be seen by comparing the three Mask R-CNN results with ResNet-50 and the two Mask R-CNN results with ResNet-101: there is a notable difference in their accuracy. Specifically, the ones reported by [75] have a much higher accuracy.

With the same ResNet-50 backbone and 1x schedule, MOSAICOS achieves $24.49/18.30$ AP/$AP_r$, better than EQL v2 [65] (23.70/14.90), BaGS [47], and cRT [38]. With the ResNet-101 backbone, MOSAICOS with 1x schedule achieves $26.54$ AP, outperforming both EQL [66] (2x

---

[12]EQL v2 [65] and Seesaw loss [75] use another implementation from [9], which uses 2x schedule for training the models on LVIS v1.0.

Table 19. **Comparisons of instance segmentation on LVIS v1.0.** MOSAICOS achieves comparable improvements against the Mask R-CNN baseline. We note that, Seesaw loss [75] uses a different implementation and training schedule (*i.e.*, 2x). Thus, the results may not be directly comparable.

| Backbone | Schedule | Method | AP | $AP_r$ | $AP_c$ | $AP_f$ |
|---|---|---|---|---|---|---|
| R-50 | 2x | Mask RCNN [24] | 23.70 | 13.50 | 22.80 | 29.30 |
| | | Seesaw [75] | (+2.70) 26.40 | (+6.10) 19.60 | (+3.30) 26.10 | (+0.50) 29.80 |
| | 1x | Mask RCNN [24] | 22.59 | 12.31 | 21.30 | 28.55 |
| | | MOSAICOS | (+1.90) 24.49 | (+5.99) 18.30 | (+1.70) 23.00 | (+0.32) 28.87 |
| R-101 | 2x | Mask RCNN [24] | 25.50 | 16.60 | 24.50 | 30.60 |
| | | Seesaw [75] | (+2.60) 28.10 | (+3.40) 20.00 | (+3.50) 28.00 | (+1.30) 31.90 |
| | 1x | Mask RCNN [24] | 24.82 | 15.18 | 23.71 | 30.31 |
| | | MOSAICOS | (+1.95) 26.77 | (+5.61) 20.79 | (+2.05) 25.76 | (+0.22) 30.53 |

schedule, 26.20 AP) and BaGS [47] (2x schedule, 25.80 AP). We also show a detailed comparison to Seesaw loss [75] in Table 19. MOSAICOS demonstrates a comparable performance gain against the Mask R-CNN baseline.

## J. Qualitative Results

We show qualitative results on LVIS v0.5 object detection in Figure 12 and Figure 13. We compare the ground truth, the results of the baseline and of our method.

We observe that our method can accurately recognize more objects from rare categories that may be overlooked by the baseline detector. For example, as shown in Figure 13, MOSAICOS correctly detects giant panda, scoreboard, horse carriage, and diaper. They are all rare classes and the baseline detector fails to make any correct detection (*i.e.*, localization and classification) on them. Moreover, the results demonstrate that MOSAICOS is able to correct the prediction labels that were wrongly classified to frequent classes without sacrificing the detection performance on common and frequent classes. As shown in the second row of Figure 12, the baseline detector wrongly predicts frequent class labels like bowl and knife with high confidence score, while MOSAICOS suppresses them and successfully predicts rare classes napkin and cappuccino.

One characteristic of LVIS is that the objects may not be exhaustively annotated in each image. We find that MOSAICOS still detects those objects which are not labeled as the ground truths. In the second and third row of Figure 13, the predictions on banner and horse are obviously correct while LVIS doesn't have annotations on them.

19

Figure 12. **Qualitative results on object detection.** Our approach can detect the rare objects missed by the baseline detector (*e.g.*, *cock*, *cappuccino*, *ferris wheel*) and correct the labels that were wrongly classified to frequent categories (*e.g.*, *bear*, *knife*, *bowl*). We superimpose **green** arrows to show where we did right while the baseline did wrong (**magenta**). **Yellow**/**Cyan**/**Red** boxes indicate frequent/common/rare (predicted) class labels.
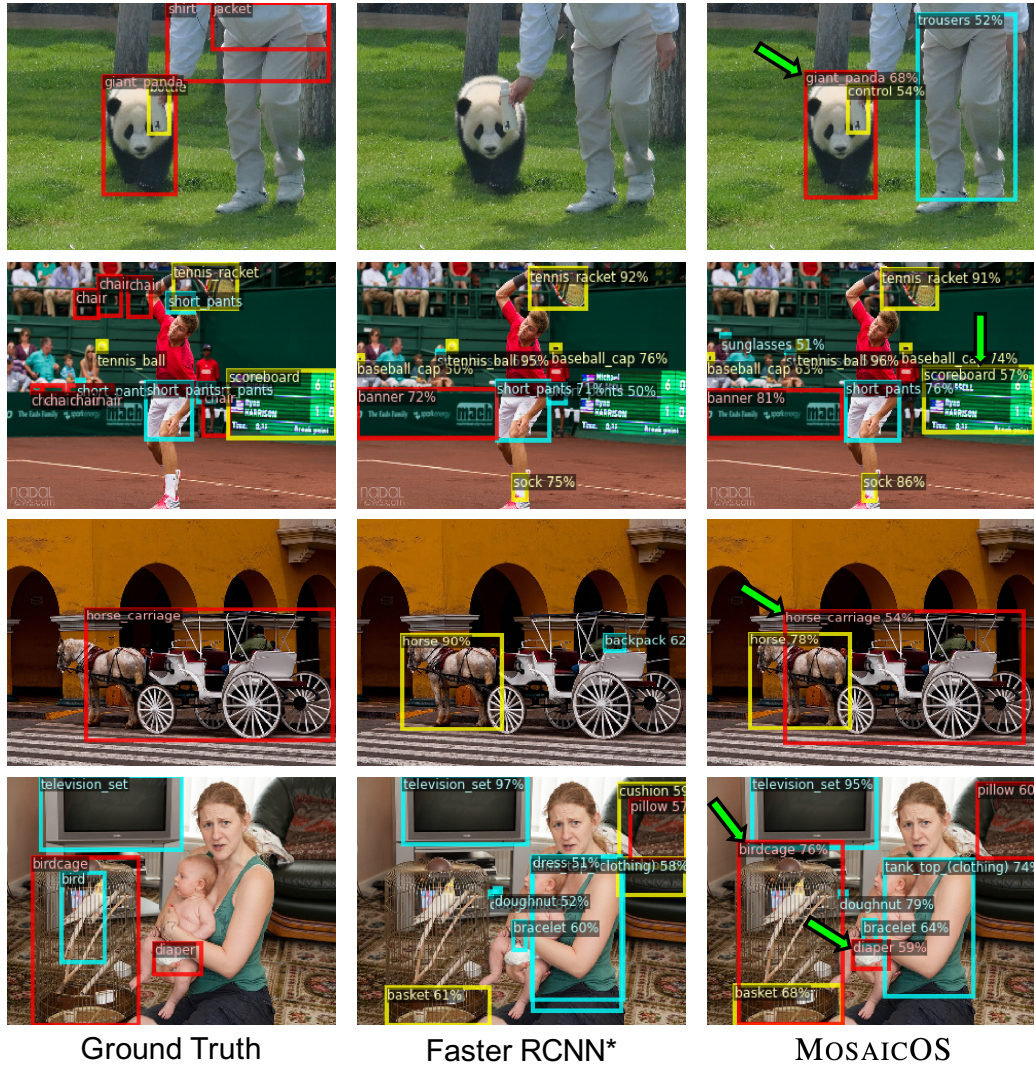
Figure 13. **Additional qualitative results on object detection.** We superimpose **green** arrows to show that our approach can detect the objects missed by the baseline detector (*e.g.*, *gaint panda*, *scoreboard*, *horse carriage*, *birdcage*, *diaper*). **Yellow**/**Cyan**/**Red** boxes indicate frequent/common/rare (predicted) class labels.