# COCONut: Modernizing COCO Segmentation
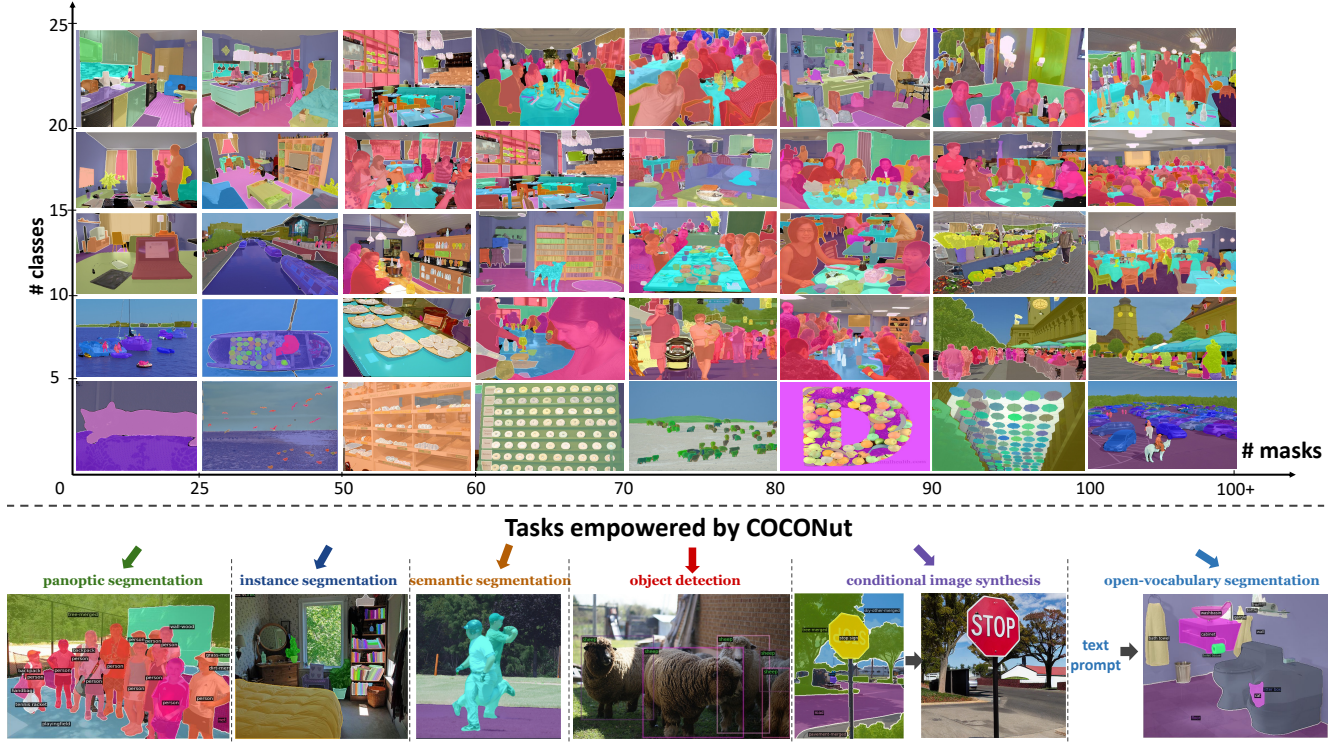
Xueqing Deng    Qihang Yu    Peng Wang    Xiaohui Shen    Liang-Chieh Chen
ByteDance
https://xdeng7.github.io/coconut.github.io/

Figure 1. **Overview of COCONut, the *COCO Next Universal segmenTation* dataset:** *Top:* COCONut, comprising images from COCO and Objects365, constitutes a diverse collection annotated with high-quality masks and semantic classes. *Bottom:* COCONut empowers a multitude of image understanding tasks.

## Abstract

*In recent decades, the vision community has witnessed remarkable progress in visual recognition, partially owing to advancements in dataset benchmarks. Notably, the established COCO benchmark has propelled the development of modern detection and segmentation systems. However, the COCO segmentation benchmark has seen comparatively slow improvement over the last decade. Originally equipped with coarse polygon annotations for 'thing' instances, it gradually incorporated coarse superpixel annotations for 'stuff' regions, which were subsequently heuristically amalgamated to yield panoptic segmentation annotations. These annotations, executed by different groups of raters, have resulted not only in coarse segmentation masks but also in inconsistencies between segmentation types. In this study, we undertake a comprehensive reevaluation of the COCO segmentation annotations. By enhancing the annotation quality and expanding the dataset to encompass 383K images with more than 5.18M panoptic masks, we introduce COCONut, the **COCO Next Universal segmenTation** dataset. COCONut harmonizes segmentation annotations across semantic, instance, and panoptic segmentation with meticulously crafted high-quality masks, and establishes a robust benchmark for all segmentation tasks. To our knowledge, COCONut stands as the inaugural large-scale universal segmentation dataset, verified by human raters. We anticipate that the release of COCONut will significantly contribute to the community's ability to assess the progress of novel neural networks.*
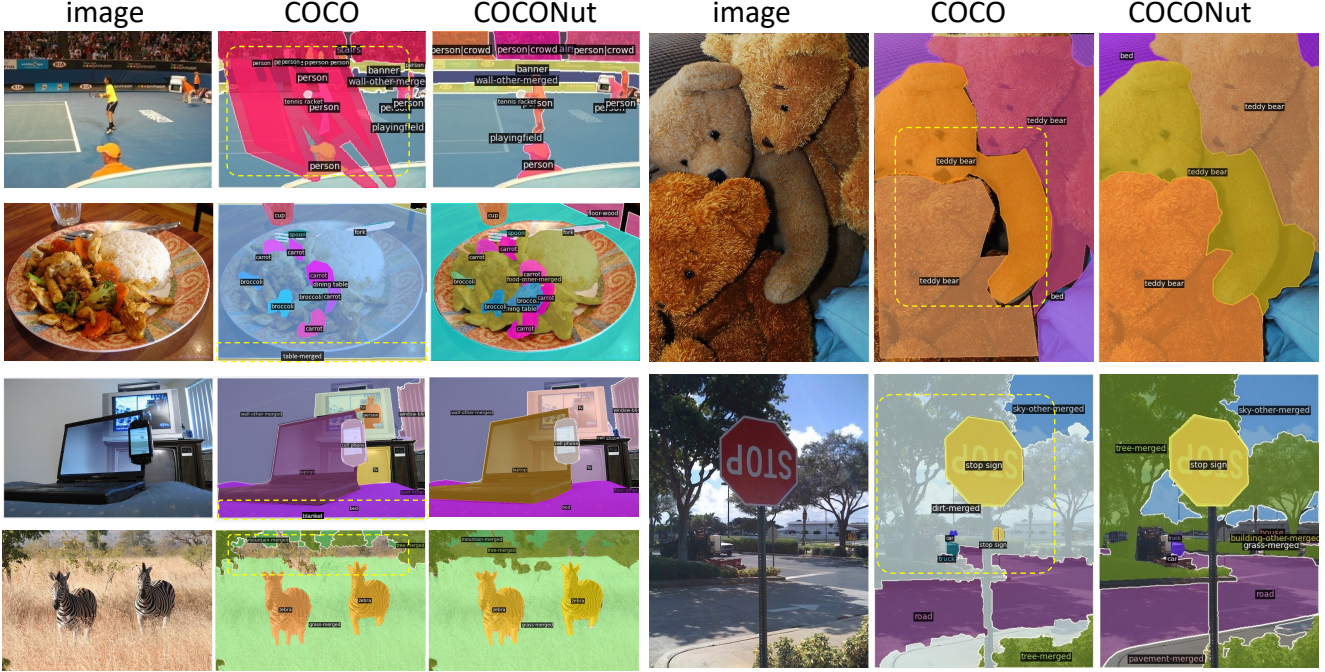
Figure 2. **Annotation Comparison:** We delineate erroneous annotations from COCO using yellow dotted line boxes, juxtaposed with our COCONut corrected annotations. Common COCO annotation errors include over-annotations (*e.g.*, 'person crowd' erroneously extends into 'playingfield'), incomplete mask fragments (*e.g.*, 'table-merged' and 'blanket' are annotated in small isolated segments), missing annotations (*e.g.*, 'tree-merged' remains unannotated), coarse segmentations (especially noticeable in 'stuff' regions annotated by superpixels and in 'thing' regions by loose polygons), and wrong semantic categories (*e.g.*, 'tree-merged' is incorrectly tagged as 'dirt-merged').

# 1. Introduction

Over the past decades, significant advancements in computer vision have been achieved, partially attributed to the establishment of comprehensive benchmark datasets. The COCO dataset [36], in particular, has played a pivotal role in the development of modern vision models, addressing a wide range of tasks such as object detection [3, 18, 22, 37, 47, 49, 71], segmentation [5–7, 10, 23, 29, 41, 55, 58–60, 64, 67], keypoint detection [20, 25, 46, 56], and image captioning [8, 48, 65]. Despite the advent of large-scale neural network models [4, 14, 40] and extensive datasets [31, 54], COCO continues to be a primary benchmark across various tasks, including image-to-text [34, 38, 65] and text-to-image [52, 66] multi-modal models. It has also been instrumental in the development of novel models, such as those fine-tuning on COCO for image captioning [50, 66] or open-vocabulary recognition [17, 19, 32, 63, 69, 70]. However, nearly a decade since its introduction, the suitability of COCO as a benchmark for contemporary models warrants reconsideration. This is particularly pertinent given the potential nuances and biases embedded within the dataset, reflective of the early stages of computer vision research.

COCO's early design inevitably encompassed certain annotation biases, including issues like imprecise object boundaries and incorrect class labels (Fig. 2). While these limitations were initially acceptable in the nascent stages of computer vision research (*e.g.*, bounding boxes are invariant to the coarsely annotated masks as long as the extreme points are the same [45]), the rapid evolution of model architectures has led to a performance plateau on the COCO benchmark[1]. This stagnation suggests a potential overfitting to the dataset's specific characteristics, raising concerns about the models' applicability to real-world data. Furthermore, despite COCO's diverse annotations supporting various tasks, its annotation is neither exhaustive nor consistent. This in-exhaustiveness is evident in the segmentation annotations, where instances of incomplete labeling are commonplace. Additionally, discrepancies between semantic, instance, and panoptic annotations within the dataset present challenges in developing a comprehensive segmentation model. Moreover, in the context of the ongoing shift towards even larger-scale datasets [16, 53], COCO's repository of approximately 120K images and 1.3M masks appears increasingly inadequate. This limitation hampers its utility in training and evaluating models designed to process and learn from substantially larger and more varied datasets.

To modernize COCO segmentation annotations, we pro-

---

[1]https://paperswithcode.com/dataset/coco

| | COCONut | COCO-17 [36] | EntitySeg [42] | ADE20K [72] | Sama-COCO [73] | LVIS [21] | Open Images [33] | COCO-Stuff [2] | PAS-21 [15] | PC-59 [43] |
|---|---|---|---|---|---|---|---|---|---|---|
| # images (train/val/test) | 358K / 25K / - | 118K / 5K / 41K | 10K / 1.5K / -† | 20K / 2K / 3K‡ | 118K / 5K / - | 100K / 20K / 40K | 944K / 13K / 40K | 118K / 5K / 41K | 1.4K / 1.4K / 1.4K | 5K / 5K / - |
| # masks / image | 13.2 / 17.4 / - | 11.2 / 11.3 / - | 16.8 / 16.4 / - | 13.4 / 15.1 / - | 9.0 / 9.5 / - | 12.7 / 12.4 / - | 2.8 / 1.8 / 1.8 | 8.6 / 8.9 / - | 2.5 / 2.5 / - | 4.9 / 4.8 / - |
| # masks | 4.75M / 437K / - | 1.3M / 57K / - | 0.17M / 24K / - | 0.27M / 30K / - | 1.07M / 47K / - | 1.27M / 0.24M / - | 2.7M / 25K / 74K | 1.02M / 44K / - | 3.6K / 3.6K / - | 24K / 24K / - |
| # thing classes | 80 | 80 | 535 | 115 | 80 | 1203 | 350 | - | - | - |
| # stuff classes | 53 | 53 | 109 | 35 | - | - | - | 91 | 21 | 59 |
| panoptic segmentation | ✓ | ✓ | ✓ | ✓ | | | | | | |
| instance segmentation | ✓ | ✓ | ✓ | △ | ✓ | ✓ | ✓ | | | |
| semantic segmentation | ✓ | △ | ✓ | ✓ | | | | ✓ | ✓ | ✓ |
| object detection | ✓ | ✓ | △ | △ | ✓ | ✓ | ✓ | | | |

Table 1. **Dataset Comparison:** We compare existing segmentation datasets that focus on daily images (street-view images are not our focus). The definition of 'thing' and 'stuff' classes are different across datasets, where the 'stuff' classes are not annotated with instance identities. †: EntitySeg dataset comprises 33K images, of which only 11K are equipped with panoptic annotations. ‡: ADE20K test server only supports semantic segmentation, and its panoptic annotations are derived by merging separately annotated instance and semantic segmentation maps, introducing minor inconsistencies between segmentation types. △: task supported, but not typically used.

pose the development of a novel, large-scale universal segmentation dataset, dubbed COCONut for the **COCO**Next **U**niversal segmen**T**ation dataset. Distinct in its approach to ensuring high-quality annotations, COCONut features human-verified mask labels for 383K images. Unlike previous attempts at creating large-scale datasets, which often compromise on label accuracy for scale [53, 61], our focus is on maintaining human verification as a standard for dataset quality. To realize this ambition, our initial step involves meticulously developing an assisted-manual annotation pipeline tailored for high-quality labeling on the subset of 118K COCO images, designated as COCONut-S split. The pipeline benefits from modern neural networks (bounding box detector [44] and mask segmenter [57, 68]), allowing our annotation raters to efficiently edit and refine those proposals. Subsequently, to expand the data size while preserving quality, we develop a data engine, leveraging the COCONut-S dataset as a high-quality training dataet to upgrade the neural networks. The process iteratively generates various sizes of COCONut training sets, yielding COCONut-B (242K images and 2.78M masks), and COCONut-L (358K images and 4.75M masks).

We adhere to a principle of consistency in annotation, aiming to establish a universal segmentation dataset (*i.e.*, consistent annotations for all panoptic/instance/semantic segmentation tasks). Additionally, COCONut includes a meticulously curated high-quality validation set, COCONut-val, comprising 5K images carefully re-labeled from the COCO validation set, along with an additional 20K images from Objects365 [54] (thus, totally 25K images and 437K masks).

To summarize, our contributions are threefold:

- We introduce COCONut, a modern, universal segmentation dataset that encompasses about 383K images and 5.18M human-verified segmentation masks. This dataset represents a significant expansion in both scale and quality of annotations compared to existing datasets. Additionally, COCONut-val, featuring meticulously curated high-quality annotations for validation, stands as a novel and challenging testbed for the research community.

- Our study includes an in-depth error analysis of the COCO dataset's annotations. This analysis not only reveals various inconsistencies and ambiguities in the existing labels but also informs our approach to refining label definitions. As a result, COCONut features ground-truth annotations with enhanced consistency and reduced label map ambiguity.

- With the COCONut dataset as our foundation, we embark on a comprehensive analysis. Our experimental results not only underscore the efficacy of scaling up datasets with high-quality annotations for both training and validation sets, but also highlight the superior value of human annotations compared to pseudo-labels.

## 2. Related Work

In this work, we focus on segmentation datasets, featuring daily images (Tab. 1). A prime example of this is the COCO dataset [36], which has been a cornerstone in computer vision for over a decade. Initially, COCO primarily focused on detection and captioning tasks [8]. Subsequent efforts have expanded its scope, refining annotations to support a wider array of tasks. For instance, COCO-Stuff [2] added semantic masks for 91 'stuff' categories, later integrated with instance masks to facilitate panoptic segmentation [30]. In addition to these expansions, several initiatives have aimed at enhancing the quality of COCO's annotations. The LVIS dataset [21] extends the number of object categories from 80 to 1,203, providing more comprehensive annotations for each image. Similarly, Sama-COCO [73] addresses the issue of low-quality masks in COCO by re-annotating instances at a finer granularity. Beyond the COCO-related datasets, there are other notable datasets contributing to diverse research scenarios, including ADE20K [72], PASCAL [15], and PASCAL-Context [43]. While these datasets have significantly advanced computer vision research, they still fall short in either annotation quality or quantity when it comes to meeting the demands for high-quality large-scale datasets.

In the realm of recent dataset innovations, SA-1B [31] stands out with its unprecedented scale, comprising 11M

3

images and 1B masks. However, a critical aspect to consider is the feasibility of human annotation at such an immense scale. Consequently, a vast majority (99.1%) of SA-1B's annotations are machine-generated and lack specific class designations. Additionally, its human annotations are not publicly released. Contrasting with scaling dataset size, the EntitySeg dataset [42] prioritizes enhancing annotation quality. This dataset features high-resolution images accompanied by meticulously curated high-quality mask annotations. However, the emphasis on the quality of annotations incurs significant resource demands, which in turn limits the dataset's scope. As a result, EntitySeg encompasses a relatively modest collection of 33K images, of which only approximately one-third are annotated with panoptic classes. Along the same direction of scaling up datasets, we present COCONut, a new large scale dataset with high quality mask annotations and semantic tags.

## 3. Constructing the COCONut Dataset

In this section, we first revisit COCO's class map definition (Sec. 3.1) and outline our image sources and varied training data sizes (Sec. 3.2). The construction of COCONut centers on two key objectives: high quality and large scale. To achieve these, we establish an efficient annotation pipeline ensuring both mask quality and accurate semantic tags (Sec. 3.3). This pipeline facilitates scalable dataset expansion while upholding annotation quality (Sec. 3.4).

### 3.1. COCO's Class Map Definition

In alignment with the COCO panoptic set [30], COCONut encompasses 133 semantic classes, with 80 categorized as 'thing' and 53 as 'stuff.' Adopting the same COCO class map ensures backward compatibility, enabling the initial use of models trained on COCO-related datasets [2, 30, 36, 73] to generate pseudo labels in our annotation pipeline.

Notably, COCONut refines class map definitions compared to COCO, offering greater clarity in our annotation instruction protocol. Building upon COCO's class map, we introduce additional definitions and instructions for labeling segmentation masks. To mitigate the annotation confusion, we meticulously define label map details and provide clear instructions to our annotation raters. For comprehensive definitions and annotation instructions for all 133 classes, please refer to the supplementary materials.

### 3.2. Image Sources and Data Splits

The images comprising COCONut are sourced from public datasets. Primarily, we aggregate images from the original COCO training and validation sets as well as its unlabeled set. Additionally, we select approximately 136K images from Objects365 dataset [54], each annotated with bounding boxes and containing at least one COCO class. This comprehensive collection results in a total of 358K and

| dataset splits | image sources | #images | #masks | #masks/image |
|---|---|---|---|---|
| COCONut-S | COCO training set [36] | 118K | 1.54M | 13.1 |
| COCONut-B | + COCO unlabeled set [36] | 242K | 2.78M | 11.5 |
| COCONut-L | + subset of Objects365 [54] | 358K | 4.75M | 13.2 |
| relabeled COCO-*val* | COCO validation set [36] | 5K | 67K | 13.4 |
| COCONut-*val* | + subset of Objects365 [54] | 25K | 437K | 17.4 |

Table 2. **Definition of COCONut Dataset Splits:** Statistics are shown accumulatively. Notably, our COCONut-val contains large #masks/image, preseting a more challenging testbed.

25K images for training and validation, respectively. As illustrated in Tab. 2, we meticulously define diverse training datasets for COCONut, spanning from 118K images to 358K images. COCONut-S (small) encompasses the same images as the original COCO training set, totaling 118K images. We adopt COCO panoptic [36] and Sama-COCO [73] masks (high-quality instance segmentation annotations[2]) as our starting point. COCONut-B (base) incorporates additional images from the COCO unlabeled set, totaling 242K images. Finally, with extra 116K images from the Objects365 dataset, COCONut-L (large) comprises 358K images. Additionally, COCONut-val contains 5K images from the COCO validation set along with an additional 20K Objects365 images.

### 3.3. Assisted-Manual Annotation Pipeline

**Annotation Challenges:** The task of densely annotating images with segmentation masks, coupled with their semantic tags (*i.e.*, classes), is exceptionally labor-intensive. Our preliminary studies reveal that, on average, it takes one expert rater approximately 5 minutes to annotate a single mask. Extrapolating this to annotate images at a scale of 10M masks would necessitate 95 years with just one expert rater. Even with a budget to employ 100 expert raters, the annotation process would still require about a year to complete. Given the extensive time and cost involved, this challenge underscores the need to explore a more effective and efficient annotation pipeline.

**Annotation Pipeline:** In response to the challenges, we introduce the assisted-manual annotation pipeline, utilizing neural networks to augment human annotators. As illustrated in Fig. 3, the pipeline encompasses four key stages: (1) machine-generated prediction, (2) human inspection and editing, (3) mask generation or refinement, and (4) quality verification. Recognizing the inherent differences between 'thing' (countable objects) and 'stuff' (amorphous regions), we meticulously address them at every stage.

**Machine-Generated Prediction:** In handling 'thing' classes, we utilize the bounding box object detector DETA [44], and for 'stuff' classes, we deploy the mask segmenter kMaX-DeepLab [68]. This stage yields a set of box proposals for 'thing' and mask proposals for 'stuff'.

---
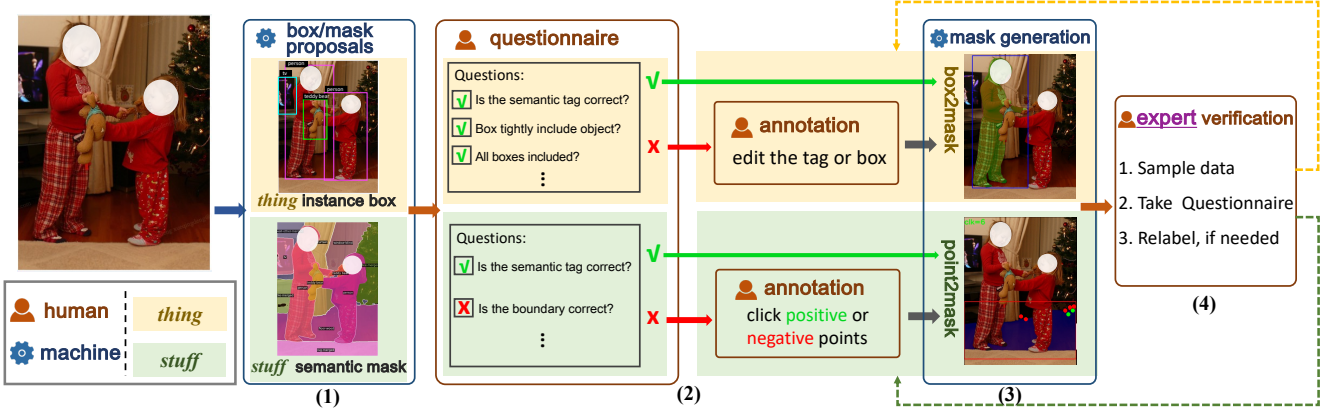
[2]https://www.sama.com/sama-coco-dataset

Figure 3. **Overview of the Proposed Assisted-Manual Annotation Pipeline:** To streamline the labor-intensive labeling task, our annotation pipeline encompasses four pivotal stages: (1) machine-generated pseudo labels, (2) human inspection and editing, (3) mask generation or refinement, and (4) quality verification. Acknowledging the inherent distinctions between 'thing' and 'stuff' classes, we systematically address these intricacies at each stage. Stage 1: Machines are employed to generate box and mask proposals for 'thing' and 'stuff', respectively. Stage 2: Raters assess the proposal qualities using a meticulously crafted questionnaire. For proposals falling short of requirements, raters can update them by editing boxes or adding positive/negative points for 'thing' and 'stuff', respectively. Stage 3: We utilize Box2Mask and Point2Mask modules to generate masks based on the inputs from stage 2. Stage 4: Experts perform a comprehensive verification of annotation quality, with relabeling done if the quality falls below our stringent standards.

**Human Inspection and Editing:** With the provided box and mask proposals, raters meticulously evaluate them based on a prepared questionnaire (*e.g.*, Is the box/mask sufficiently accurate? Is the tag correct? Any missing boxes?) The raters adhere to stringent standards during inspection to ensure proposal quality. In cases where proposals fall short, raters are directed to perform further editing. Specifically, for 'thing' classes, raters have the flexibility to add or remove boxes along with their corresponding tags (i.e., classes). In the case of 'stuff' classes, raters can refine masks by clicking positive or negative points, indicating whether the points belong to the target instance or not.

**Mask Generation or Refinement:** Utilizing the provided boxes and masks from the preceding stage, we employ the **Box2Mask** and **Point2Mask** modules to generate segmentation masks for 'thing' and 'stuff' classes, respectively. The **Box2Mask** module extends kMaX-DeepLab, resulting in the box-kMaX model, which generates masks based on provided bounding boxes. This model incorporates additional box queries in conjunction with the original object queries. The added box queries function similarly to the original object queries, except that they are initialized using features pooled from the backbone within the box regions (original object queries are randomly initialized). As shown in Fig. 4, leveraging object-aware box queries enables box-kMaX to effectively segment 'thing' objects with the provided bounding boxes. The **Point2Mask** module utilizes the interactive segmenter CFR [57], taking positive/negative points as input and optionally any initial mask (from either kMaX-DeepLab or the previous round's output
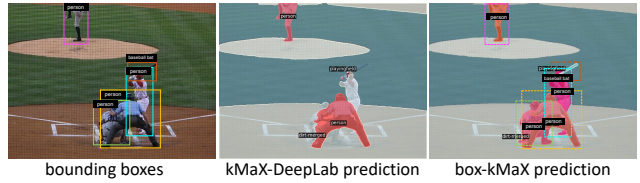


Figure 4. **Mask Prediction Comparison:** In contrast to kMaX-DeepLab, box-kMaX (Box2Mask module) leverages box queries, initialized with features pooled from the backbone within the box regions, enabling more accurate segmentation of 'thing' objects. Notably, kMaX-DeepLab falls short in capturing the challenging 'baseball bat' and the heavily occluded 'person' in the figure.

mask). This stage allows us to amass a collection of masks generated from boxes and refined by points.

It is worth noting that there are other interactive segmenters that are also capable of generating masks using box and point as inputs (*e.g.*, SAM [31], SAM-HQ[28]). However, our analyses (in Sec. 4) indicate that the tools we have developed suffice for our raters to produce high-quality annotations. The primary focus of our work is to conduct a comprehensive analysis between the original COCO dataset and our newly annotated COCONut. Improving interactive segmenters lies outside the scope of this study.

**Quality Verification by Experts:** Armed with the amassed masks from the preceding stage, we task *expert raters* with quality verification. Unlike the general human raters in stage 2, our expert raters boast extensive experience in dense pixel labeling (5 years of proficiency in
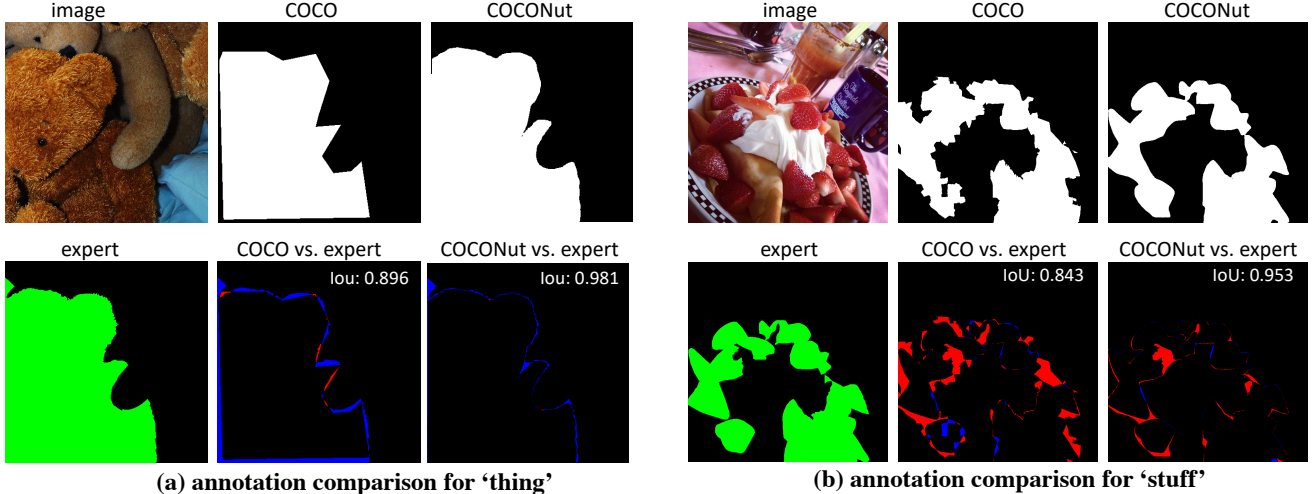
(a) annotation comparison for 'thing'  (b) annotation comparison for 'stuff'

Figure 5. **Annotation Comparison:** We show annotations obtained by COCO, COCONut (Box2Mask for 'thing' in (a) or Point2Mask for 'stuff' in (b)), and our expert rater. COCONut's annotation exhibits sharper boundaries, closely resembling expert results, as evident from higher IoU values. The blue and red regions correspond to extra and missing regions, respectively, compared to the expert mask.

Photoshop). To manage the extensive volume of annotated masks with only two experts, we opt for a random sampling of 50%. The experts meticulously assess these masks, along with their associated tags, using the same carefully crafted questionnaire as in the previous stage. Furthermore, recognizing the Box2Mask module's reliance on provided bounding boxes, we additionally instruct experts to verify the accuracy of box proposals, selecting 30% samples for a thorough quality check. Should any fall short of our stringent requirements, they undergo relabeling using the time-intensive Photoshop tool to ensure high annotation quality.

### 3.4. Data Engine for Scaling Up Dataset Size

**Overview:** With the streamlined assisted-manual annotation pipeline in place, we build a data engine to facilitate the dataset expansion. Our data engine capitalizes on the annotation pipeline to accumulate extensive, high-quality annotations, subsequently enhancing the training of new neural networks for improved pseudo-label generation. This positive feedback loop is iteratively applied multiple times.

**Data Engine:** Machines play a crucial role in generating box/mask proposals (stage 1) and refined masks (stage 3) in the assisted-manual annotation pipeline. Initially, publicly available pre-trained neural networks are employed to produce proposals. Specifically, DETA [44] (utilizing a Swin-L backbone [39] trained with Objects365 [54] and COCO detection set [36]) and kMaX-DeepLab [68] (featuring a ConvNeXt-L backbone [40] trained with COCO panoptic set [30]) are utilized to generate box and mask proposals for 'thing' and 'stuff', respectively. The Point2Mask module (built upon CFR [57]) remains fixed throughout the COCONut construction, while the Box2Mask module (box-

|  | 'thing' | 'stuff' |
|---|---|---|
| expert-1 vs. expert-2 | 98.1% | 97.3% |
| raters vs. experts | 96.3% | 96.7% |

(a) **Annotation Agreement**

|  | 'thing' | 'stuff' |
|---|---|---|
| purely-manual | 10 min | 5 min |
| assisted-manual | 10 sec | 42 sec |

(b) **Annotation Speed**

Table 3. **Annotation Analysis:** (a) Our two experts and raters demonstrate a high level of agreement in their annotations. (b) The assisted-manual pipeline expedites the annotation.

kMaX, a variant of kMaX-DeepLab using box queries) is trained on COCO panoptic set. The annotation pipeline initially produces the COCONut-S dataset split. Subsequently, COCONut-S is used to re-train kMaX-DeepLab and box-kMaX, enhancing mask proposals for 'stuff' and Box2Mask capabilities, respectively. Notably, DETA and the Point2Mask module are not re-trained, as DETA is already pre-trained on a substantial dataset, and CFR exhibits robust generalizability. The upgraded neural networks yield improved proposals and mask generations, enhancing the assisted-manual annotation pipeline and leading to the creation of COCONut-B. This process is iterated to generate the final COCONut-L, which also benefits from the ground-truth boxes provided by Objects365.

## 4. Annotation and Data Engine Analysis

In this section, we scrutinize the annotations produced through our proposed assisted-manual annotation pipeline (Sec. 4.1). Subsequently, we delve into the analysis of the improvement brought by our data engine (Sec. 4.2).
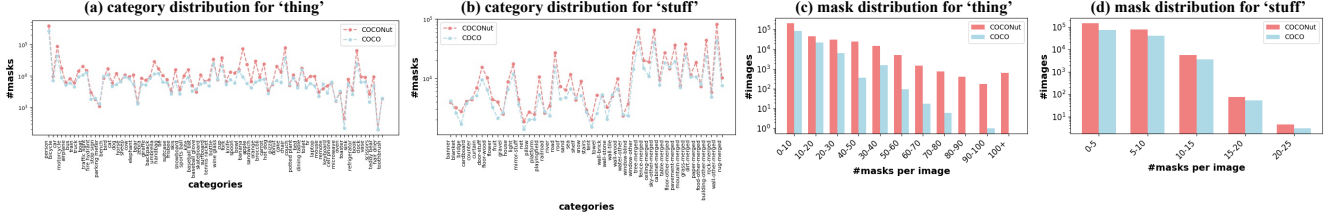
Figure 6. **Dataset Statistics:** In subfigures (a) and (b), depicting category distributions for 'thing' and 'stuff', COCONut consistently displays a higher number of masks across all categories compared to COCO. Subfigures (c) and (d) show mask distribution for 'thing' and 'stuff', respectively, demonstrating that COCONut contains a greater number of images with a higher density of masks per image.

| constructed dataset | mean | median |
|---|---|---|
| COCONut-S | 78% | 75% |
| COCONut-B | 51% | 55% |
| COCONut-L | 43% | 45% |

(a) **Non-Pass Rate in Stage 2**

| constructed dataset | mean | median |
|---|---|---|
| COCONut-S | 2.4 | 2 |
| COCONut-B | 0.8 | 1 |
| COCONut-L | 0.5 | 1 |

(b) **#Rounds of Relabeling in Stage 4**

Table 4. **Data Engine Analysis:** During the creation of the current dataset split, the mask proposals stem from models trained on datasets from preceding stages, such as COCONut-S utilizing proposal models from COCO, and so forth.

## 4.1. Annotation Analysis

**Assisted-Manual *vs*. Purely-Manual:** We conduct a thorough comparison in this study between annotations generated by our assisted-manual and purely-manual annotation pipelines. Our assessment is based on two metrics: annotation quality and processing speed.

The purely-manual annotation pipeline involves two in-house experts, each with over 5 years of experience using Photoshop for labeling dense segmentation maps. They received detailed instructions based on our annotation guidelines and subsequently served as tutorial training mentors for our annotation raters. Additionally, they played a crucial role in the quality verification of masks during stage 4.

To conduct the "agreement" experiments, we randomly selected 1000 segmentation masks and tasked our two in-house experts with annotating each mask. An "agreement" was achieved when both annotations exhibited an IoU (Intersection-over-Union) greater than 95%. As presented in Tab. 3a, our experts consistently demonstrated a high level of agreement in annotating both 'thing' and 'stuff' masks. Comparatively, minor disparities were observed in the annotations provided by our raters, highlighting their proficiency. Additionally, Tab. 3b showcases the annotation speed. The assisted-manual pipeline notably accelerates the annotation process by editing boxes and points, particularly beneficial for 'thing' annotations. Annotating 'stuff', however, involves additional time due to revising the coarse superpixel annotations by COCO. Finally, Fig. 5 presents annotation examples from COCO, our experts, and COCONut (our raters with the assisted-manual pipeline),
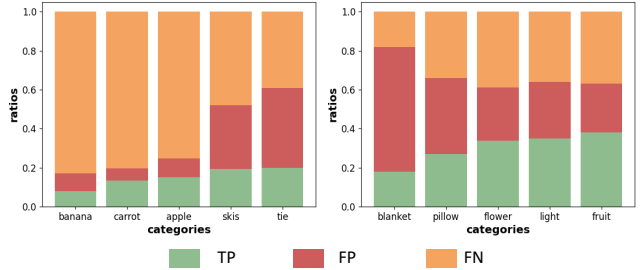


Figure 7. **Top 5 Disagreed Categories Between COCO-val and relabeled COCO-val:** COCO-val is treated as the prediction, while relabeled COCO-val serves as ground truth. The comparison showcases True Positive (TP), False Positive (FP), and False Negative (FN) rates for both 'thing' (left) and 'stuff' (right).

underscoring the high-quality masks produced.

## 4.2. Data Engine Analysis

The data engine enhances neural networks using annotated high-quality data, resulting in improved pseudo masks and decreased workload for human raters. To measure its impact, we present non-pass rates in stage 2 human inspection. These rates indicate the percentage of machine-generated proposals that failed our questionnaire's standards and required further editing. Tab. 4a demonstrates that including more high-quality training data improves non-pass rates, signifying enhanced proposal quality. Furthermore, Tab. 4b showcases the number of relabeling rounds in stage 4 expert verification, reflecting additional iterations required for annotations failing expert verification. Consistently, we observed reduced relabeling rounds with increased inclusion of high-quality training data.

## 5. Dataset Statistics

**Class and Mask Distribution:** Fig. 6 depicts the category and mask distribution within COCONut. Panels (a) and (b) demonstrate that COCONut surpasses COCO in the number of masks across all categories. Additionally, panels (c) and (d) feature histograms depicting the frequency of 'masks per image'. These histograms highlight a notable

| | PQ | SQ | RQ | $PQ^{bdry}$ | $SQ^{bdry}$ | $RQ^{bdr}$ |
|---|---|---|---|---|---|---|
| all | 67.1 | 86.2 | 77.4 | 59.2 | 79.4 | 74.5 |
| thing | 65.0 | 86.0 | 75.2 | 58.6 | 80.7 | 72.4 |
| stuff | 70.2 | 86.5 | 80.8 | 60.1 | 77.3 | 77.6 |

Table 5. **Quantitative Comparison Between COCO-val and relabeled COCO-val:** COCO-val serves as the prediction, contrasting with relabeled COCO-val as the ground-truth.
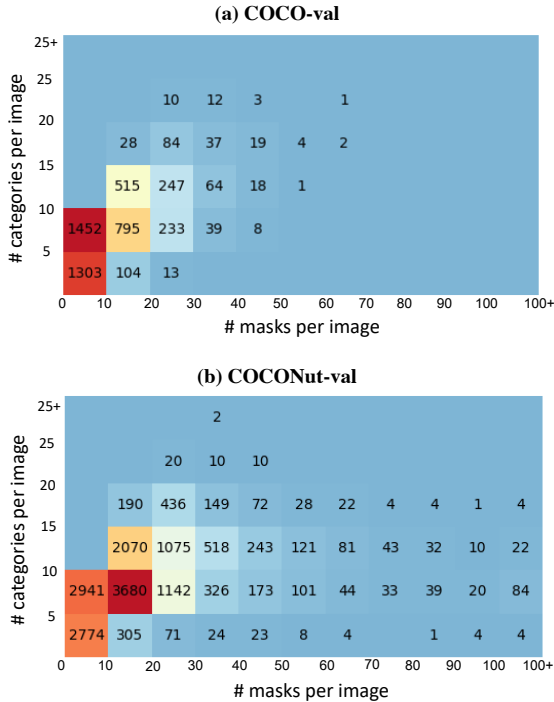


(a) COCO-val

(b) COCONut-val

Figure 8. **Mask and Class Frequency Distribution:** COCONut-val introduces a more challenging testbed compared to the original COCO-val. It features a greater number of images that contain higher quantities of both masks and distinct categories per image.

| backbone | training set | COCO-val | | | relabeled COCO-val | | | COCONut-val | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | PQ | $AP^{mask}$ | mIoU | PQ | $AP^{mask}$ | mIoU | PQ | $AP^{mask}$ | mIoU |
| ResNet50 | COCO | 53.3 | 39.6 | 61.7 | 55.1 | 40.6 | 63.9 | 53.1 | 37.1 | 62.5 |
| | COCONut-S | 51.7 | 37.5 | 59.4 | 58.9 | 44.4 | 64.4 | 56.7 | 41.2 | 63.6 |
| | COCONut-B | 53.4 | 39.3 | 62.6 | 60.2 | 45.2 | 65.7 | 58.1 | 42.9 | 64.7 |
| | COCONut-L | 54.1 | 40.2 | 63.1 | 60.7 | 45.8 | 66.1 | 60.7 | 44.8 | 68.3 |
| ConvNeXt-L | COCO | 57.9 | 45.0 | 66.9 | 60.4 | 46.4 | 69.9 | 58.3 | 44.1 | 66.4 |
| | COCONut-S | 55.9 | 41.9 | 66.1 | 64.4 | 50.8 | 71.4 | 59.4 | 45.7 | 67.8 |
| | COCONut-B | 57.8 | 44.8 | 66.6 | 64.9 | 51.2 | 71.8 | 61.3 | 46.5 | 69.5 |
| | COCONut-L | 58.1 | 45.3 | 67.3 | 65.1 | 51.4 | 71.9 | 62.7 | 47.6 | 70.6 |

Table 6. **Training Data and Backbones:** The evaluations are conducted on three different validation sets: original COCO-val, relabeled COCO-val (by our raters), and COCONut-val.

cating numerous small isolated masks, echoing our earlier findings regarding 'bed' and 'blanket' conflicts, as depicted in Fig. 2 (row 3). Finally, Tab. 5 provides a quantitative analysis comparing COCO-val and our relabeled COCO-val. The results emphasize the notable divergence between the two sets, underscoring our dedicated efforts to improve the annotation quality of validation set. The discrepancy is particularly evident in boundary metrics [11]. Notably, the divergence in stuff $SQ^{bdry}$ reflects our enhancements to the original 'stuff' annotations by superpixels [1, 2].

**COCONut-val (a new challenging testbed):** To augment our relabeled COCO-val, we introduced an additional 20K annotated images from Objects365, forming COCONut-val. Fig. 8 illustrates 2D histograms comparing COCO-val and COCONut-val, where we count the number of images w.r.t. their #masks and #categories per image. The figure showcases that COCO-val annotations are concentrated around a smaller number of masks and categories, whereas COCONut-val demonstrates a broader distribution, with more images having over 30 masks. On average, COCONut-val boasts 17.4 masks per image, significantly exceeding COCO-val's average of 11.3 masks.

# 6. Discussion

In light of the COCONut dataset, we undertake a meticulous analysis to address the following inquiries. We employ kMaX-DeepLab [68] throughout the experiments, benchmarked with several training and validation sets.

**COCO encompasses only 133 semantic classes. Is an extensive collection of human annotations truly necessary?** We approach this query from two vantage points: the training and validation sets. Tab. 6 showcases consistent improvements across various backbones (ResNet50 [24] and ConvNeXt-L [40]) and three evaluated validation sets (measured in PQ, AP, and mIoU) as the training set size increases from COCONut-S to COCONut-L. Interestingly, relying solely on the original small-scale COCO training set yields unsatisfactory performance on both relabeled COCO-val and COCONut-val sets, emphasizing the need for more human annotations in training. Despite annota-

trend in COCONut, indicating a higher prevalence of images with denser mask annotations compared to COCO.

**COCO-val *vs*. relabeled COCO-val:** We conducted a comparative analysis between the original COCO-val annotations and our relabeled COCO-val. Exploiting the Panoptic Quality (PQ) metric, we employed its True Positive (TP), False Positive (FP), and False Negative (FN) rates to assess each category. TP signifies agreement between annotations, while FP and FN highlight additional or missing masks, respectively. In Fig. 7, we present the top 5 categories displaying discrepancies for both 'thing' and 'stuff'. All these categories exhibit notably low TP rates, indicating substantial differences between COCO-val and our relabeled version. In 'thing' categories, high FN rates (around 0.8) are observed for 'banana', 'carrot', and 'apple', suggesting numerous missing masks. Conversely, 'stuff' categories exhibit high FP rates for 'blanket' and 'pillow', indi-

| backbone | training set | COCO-val | | | relabeled COCO-val | | | COCONut-val | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | PQ | AP$^{mask}$ | mIoU | PQ | AP$^{mask}$ | mIoU | PQ | AP$^{mask}$ | mIoU |
| ConvNeXt-L | COCO | 57.9 | 45.0 | 66.9 | 60.4 | 46.4 | 69.9 | 58.3 | 44.1 | 66.4 |
| | COCO-B$_M$ | 58.0 | 44.9 | 67.1 | 60.7 | 46.3 | 70.5 | 58.5 | 44.2 | 66.4 |
| | COCONut-S | 55.9 | 41.9 | 66.1 | 64.4 | 50.8 | 71.4 | 59.4 | 45.7 | 67.8 |
| | COCONut-B$_M$ | 56.2 | 41.8 | 66.3 | 64.5 | 50.9 | 71.4 | 59.5 | 45.3 | 67.7 |
| | COCONut-B | 57.8 | 44.8 | 66.6 | 64.9 | 51.2 | 71.8 | 61.3 | 46.5 | 69.5 |

Table 7. **Pseudo-Labels *vs*. Human Labels:** COCO-B$_M$ comprises the original COCO training set plus the machine pseudo-labeled COCO unlabeled set. COCONut-B$_M$ contains COCONut-S and machine pseudo-labeled COCO unlabeled set (in contrast to the fully human-labeled COCONut-B).

tion biases between COCO and COCONut (Fig. 9), training with COCONut-B achieves performance akin to the original COCO training set on the COCO validation set, hinting that a larger training corpus might mitigate inter-dataset biases.

Shifting our focus from the training set to the validation set, the results in Tab. 6 indicate performance saturation on both COCO-val and relabeled COCO-val as the training set expands from COCONut-B to COCONut-L. This saturation phenomenon in COCO-val, consisting of only 5K images, is also observed in the literature[3], suggesting its inadequacy in evaluating modern segmenters. Conversely, the newly introduced COCONut-val, comprising 25K images with denser mask annotations, significantly improves benchmarking for models trained with varied data amounts. This outcome underscores the significance of incorporating more human-annotated, challenging validation images for robust model assessment. Therefore, the inclusion of additional human-annotated images is pivotal for both training and validation, significantly impacting the performance of modern segmentation models.

**Are pseudo-labels a cost-effective alternative to human annotations?** While expanding datasets using machine-generated pseudo-labels seems promising for scaling models trained on large-scale data, its effectiveness remains uncertain. To address this, we conducted experiments outlined in Tab. 7. Initially, leveraging a checkpoint (row 1: 57.9% PQ on COCO-val), we generated pseudo-labels for the COCO unlabeled set, augmenting the original COCO training set to create the COCO-B$_M$ dataset. Surprisingly, training on COCO-B$_M$ resulted in only a marginal 0.1% PQ improvement on COCO-val, consistent across all tested validation sets (1st and 2nd rows in the table).

We hypothesized that the annotation quality of the pretrained dataset might influence pseudo-label quality. To investigate, we then utilized a different checkpoint (row 3: 64.4% PQ on relabeled COCO-val) to generate new pseudo-labels for the COCO unlabeled set. Combining these with COCONut-S produced the COCONut-B$_M$ dataset, yet still yielded a mere 0.1% PQ improvement on the relabeled

---



COCO-trained prediction
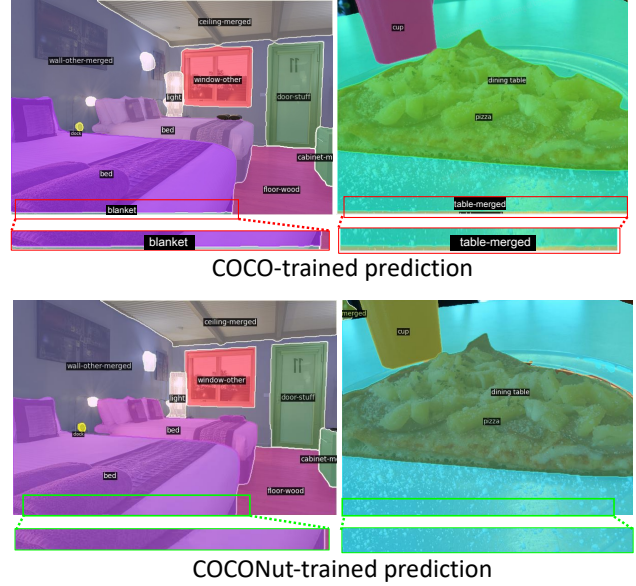


COCONut-trained prediction

Figure 9. **Influence of Training Data on Predictions:** We present predictions from two models: one trained on original COCO (top) and the other on COCONut (bottom). *Top*: The COCO-trained model predicts a small isolated mask, influenced by the biases inherent in the COCO coarse annotations (*e.g.*, see Fig. 2, row 3). *Bottom*: The COCONut-trained model does not predict small isolated masks, thanks to the meticulously crafted annotations. Best zoomed-in.

COCO-val. Notably, employing the fully human-labeled COCONut-B resulted in the most significant improvements (last row in the table). Our findings suggest limited benefits from incorporating pseudo-labels. Training with pseudo-labels seems akin to distilling knowledge from a pre-trained network [27], offering minimal additional information for training new models.

# 7. Visualization of COCONut Annotations

We present annotation visualizations for COCONut dataset. Specifically, Fig. 10 and Fig. 11 demonstrate the COCONut annotations for images sourced from COCO unlabeled set [36] and Objects365 [54]. As shown in the figures, COCONut provides annotations comprising a large number of classes and masks. Notably, the inclusion of Objects365 images enriches COCONut annotations by introducing a wider variety of classes and masks compared to the COCO images. Finally, Fig. 12 compares the COCO and COCONut annotations, where the common errors of COCO (*e.g.*, inaccurate boundary, loose polygon, missing masks, and wrong semantics) are all corrected in COCONut annotations.

---

Figure 10. **Visualization of COCONut Annotations:** This figure demonstrates COCONut annotations with images sourced from COCO unlabeled set images. COCONut provides annotations comprising a large number of classes and masks.
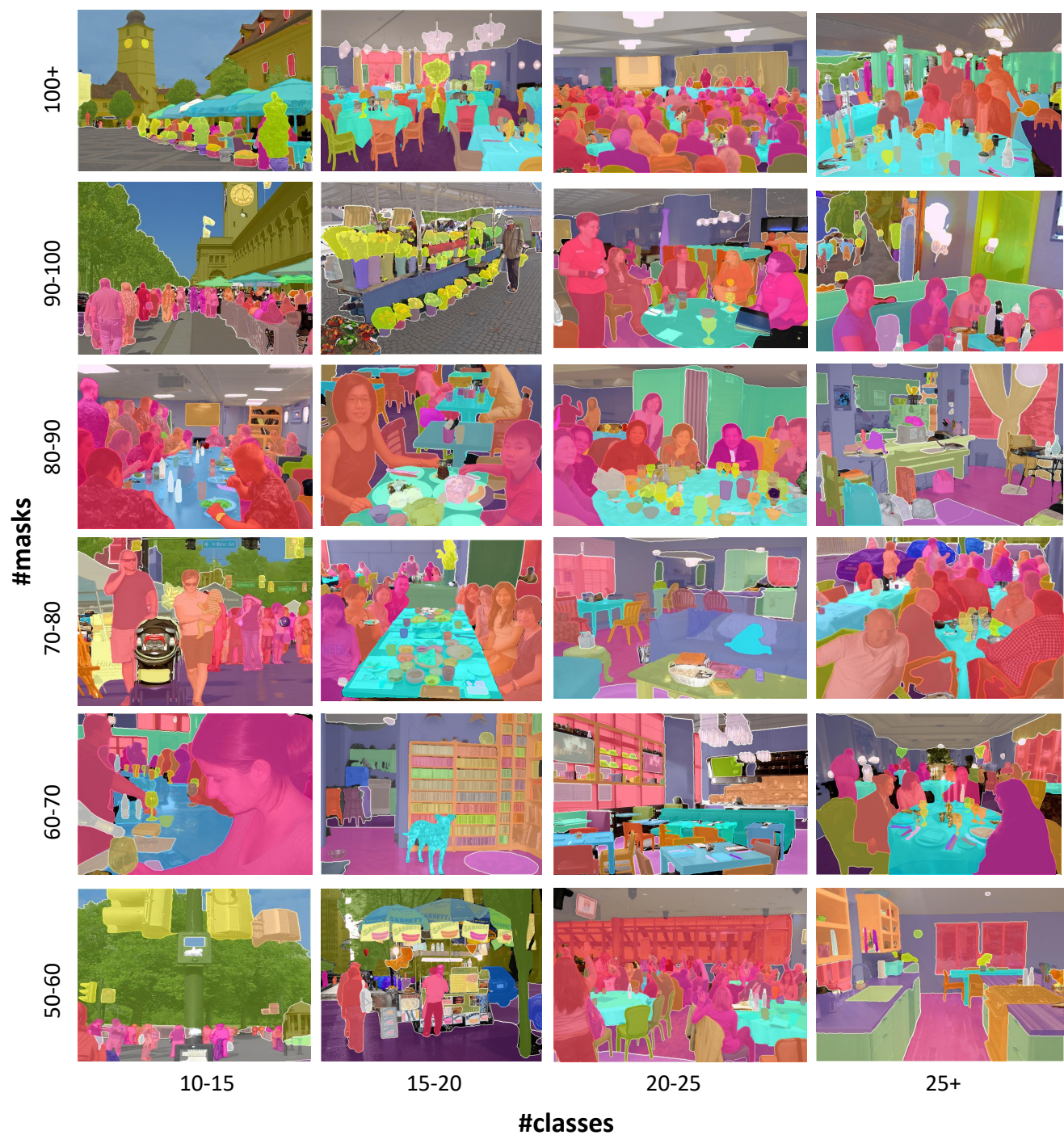
Figure 11. **Visualization of COCONut Annotations:** This figure showcases COCONut annotations using images sourced from both the COCO unlabeled set and Objects365 images. The inclusion of Objects365 images enriches COCONut annotations by introducing a wider variety of classes and masks compared to the COCO images.
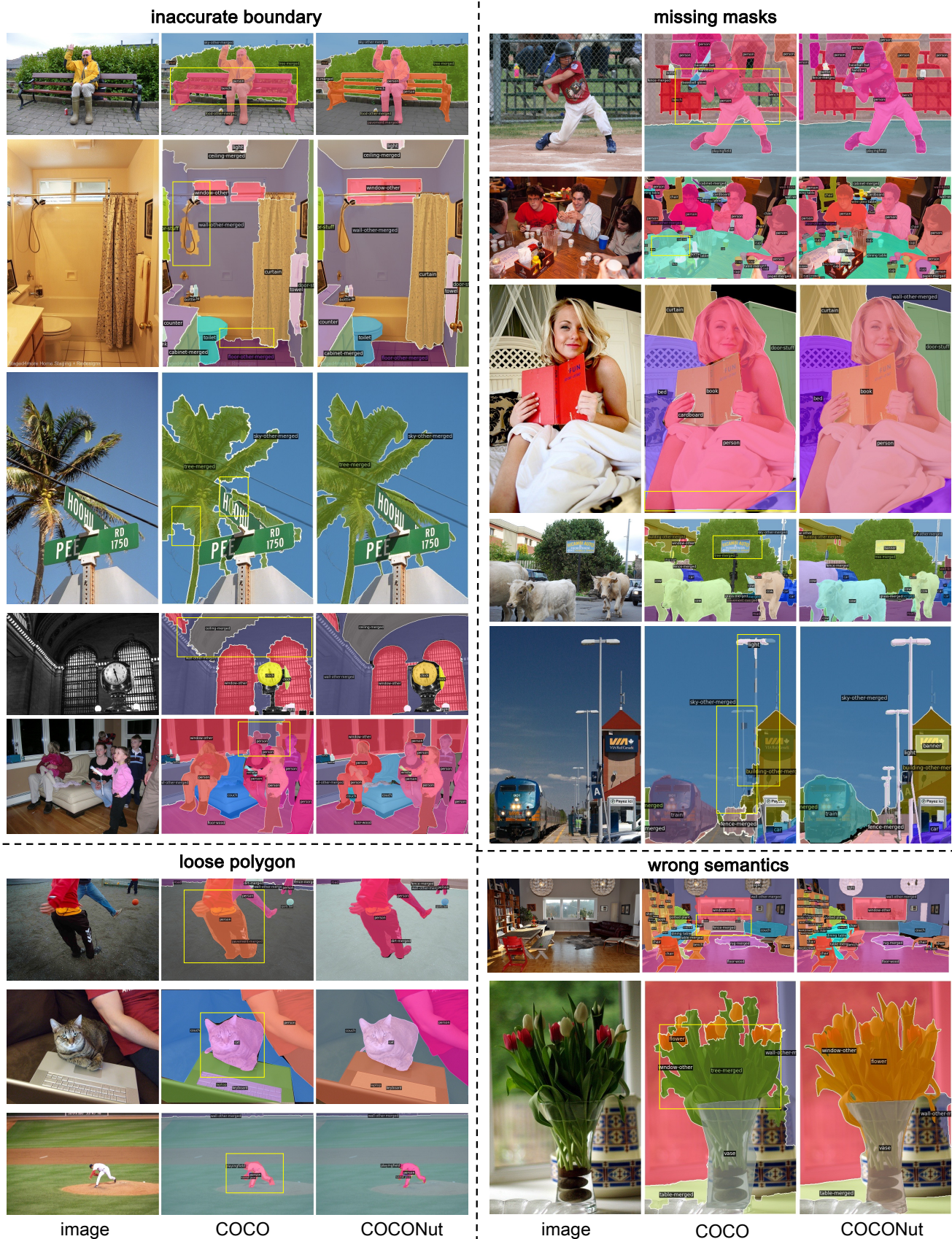
Figure 12. **Visualization Comparison Between COCO and COCONut:** COCONut effectively mitigates the annotations errors by COCO. The yellow boxes highlight the erroneous areas in COCO.

# Appendix

In the supplementary materials, we provide additional information, as listed below.

- Sec. A presents additional annotation visualizations.
- Sec. B provides more benchmarking results on CO-CONut.
- Sec. C provides the details of our label map definition and annotation rules.

## A. Additional Annotation Visualizations

We present additional annotation visualizations for CO-CONut dataset.

In particular, Fig. 13 and Fig. 14 provide more annotation comparisons between COCO, COCONut, and our expert raters. Fig. 15 and Fig. 16 provide more visualizations of prediction bias introduced by training data.

## B. Additional Experimental Results

In this section, we outline the training and evaluation protocols utilized to benchmark COCONut across multiple tasks and the corresponding results in Sec. B.1 and Sec. B.2, respectively.

### B.1. Training and Evaluation Protocols

**Training Protocol:** COCONut undertakes benchmarking across various tasks, encompassing panoptic segmentation [30], instance segmentation [22], semantic segmentation [26], object detection [15], and open-vocabulary segmentation [13, 17]. The kMaX-DeepLab [62, 68], tailored for universal segmentation, serves as the primary framework for panoptic, instance, and semantic segmentation in our experiments. Object detection relies on the DETA framework [44], while open-vocabulary segmentation utilizes the FC-CLIP framework [69].

Throughout each task, we strictly adhere to the training hyper-parameters defined by the respective frameworks, utilizing ResNet50 [24], Swin-L [39], and ConvNeXt-L [40] as the backbones.

**Evaluation Protocol:** When evaluating each task, we follow official settings meticulously. For panoptic, instance, and semantic segmentation tasks, metrics such as panoptic quality (PQ) [30], $AP^{mask}$ [36], and mean Intersection-over-Union (mIoU) [15] are reported. Bounding box detection performance is measured using the $AP^{box}$ metric. In line with prior methodologies [13, 63], open-vocabulary segmentation results undertake zero-shot evaluation on other segmentation datasets [15, 43, 72].

### B.2. COCONut Empowers Various Tasks

In this section, we show the results for task-specific models trained on COCONut datasets including panoptic segmenta-

| method | backbone | training set | COCO-val PQ | relabeled COCO-val PQ | COCONut-val PQ |
|---|---|---|---|---|---|
| kMaX-DeepLab | ResNet50 | COCO | 53.3 | 55.1 | 53.1 |
| | | COCONut-S | 51.7 | 58.9 | 56.7 |
| | | COCONut-B | 53.4 | 60.2 | 58.1 |
| | | COCONut-L | 54.1 | 60.7 | 60.7 |
| | ConvNeXt-L | COCO | 57.9 | 60.4 | 58.3 |
| | | COCONut-S | 55.9 | 64.4 | 59.4 |
| | | COCONut-B | 57.8 | 64.9 | 61.3 |
| | | COCONut-L | 58.1 | 65.1 | 62.7 |

Table 8. **Benchmarking Task-Specific Panoptic Segmentation Models:** kMaX-DeepLab is trained with *panoptic* segmentation annotations across various training and validation sets.

| method | backbone | training set | COCO-val $AP^{mask}$ | relabeled COCO-val $AP^{mask}$ | COCONut-val $AP^{mask}$ |
|---|---|---|---|---|---|
| kMaX-DeepLab | ResNet50 | COCO | 44.1 | 44.6 | 41.9 |
| | | COCONut-S | 40.9 | 49.2 | 44.9 |
| | | COCONut-B | 41.2 | 50.3 | 46.2 |
| | | COCONut-L | 41.4 | 50.9 | 47.1 |
| | ConvNeXt-L | COCO | 49.2 | 50.2 | 47.1 |
| | | COCONut-S | 45.5 | 55.8 | 51.2 |
| | | COCONut-B | 46.4 | 56.7 | 52.9 |
| | | COCONut-L | 47.0 | 57.0 | 53.8 |

Table 9. **Benchmarking Task-Specific Instance Segmentation Models:** kMaX-DeepLab is trained with *instance* segmentation annotations across various training and validation sets.

| method | backbone | training set | COCO-val mIoU | relabeled COCO-val mIoU | COCONut-val mIoU |
|---|---|---|---|---|---|
| kMaX-DeepLab | ResNet50 | COCO | 59.5 | 64.6 | 62.9 |
| | | COCONut-S | 59.3 | 66.4 | 65.1 |
| | | COCONut-B | 63.5 | 67.3 | 66.5 |
| | | COCONut-L | 64.2 | 68.0 | 67.8 |
| | ConvNeXt-L | COCO | 67.1 | 70.9 | 68.1 |
| | | COCONut-S | 66.1 | 71.9 | 69.9 |
| | | COCONut-B | 67.4 | 72.4 | 71.3 |
| | | COCONut-L | 67.5 | 72.7 | 72.6 |

Table 10. **Benchmarking Task-Specific Semantic Segmentation Models:** kMaX-DeepLab is trained with *semantic* segmentation annotations across various training and validation sets.

| backbone | training dataset | evaluation set (mIoU) | | |
|---|---|---|---|---|
| | | COCO-val | relabeled COCO-val | COCONut-val |
| ViT-Adapter-B | COCO | 61.2 | 64.5 | 61.8 |
| | COCONut-S | 60.6 | 66.0 | 64.9 |
| | COCONut-B | 61.3 | 66.9 | 66.3 |
| | COCONut-L | 62.4 | 67.7 | 67.1 |
| ViT-Adapter-L | COCO | 66.6 | 69.9 | 67.5 |
| | COCONut-S | 65.2 | 71.0 | 69.5 |
| | COCONut-B | 66.4 | 72.1 | 70.7 |
| | COCONut-L | 67.2 | 72.3 | 71.0 |

Table 11. **Benchmarking plain ViT backbone for Semantic Segmentation:** Mask2Former w/ ViT-Adapter is trained with *semantic* segmentation annotations.

tion, instance segmentation, semantic segmentation, object detection, semantic mask conditional image synthesis.

**Panoptic Segmentation:** In Tab. 8, we benchmark kMaX-DeepLab on the task of panoptic segmentation. The results are the same as Tab. 6 in the main paper, where a panoptic model is evaluated on all three segmentation metrics.

**Instance Segmentation:** We benchmark kMaX-DeepLab on the task of instance segmentation. Differ-
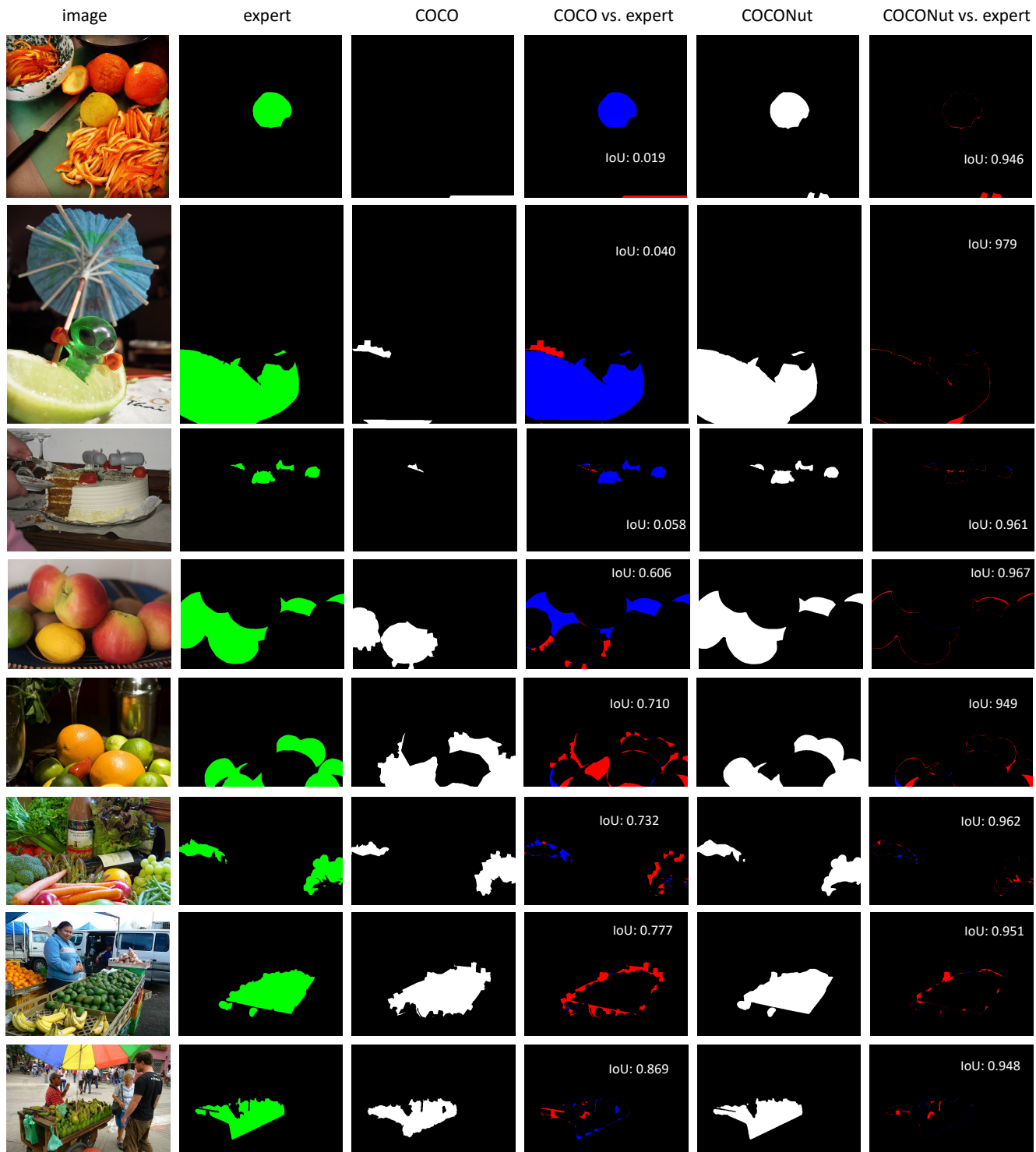
Figure 13. **Annotation Comparison:** We show annotations obtained by COCO, COCONut with Point2Mask for 'stuff', and our expert rater. COCONut's annotation exhibits sharper boundaries, closely resembling expert results, as evident from higher IoU values. The blue and red regions correspond to extra and missing regions, respectively, compared to the expert mask.
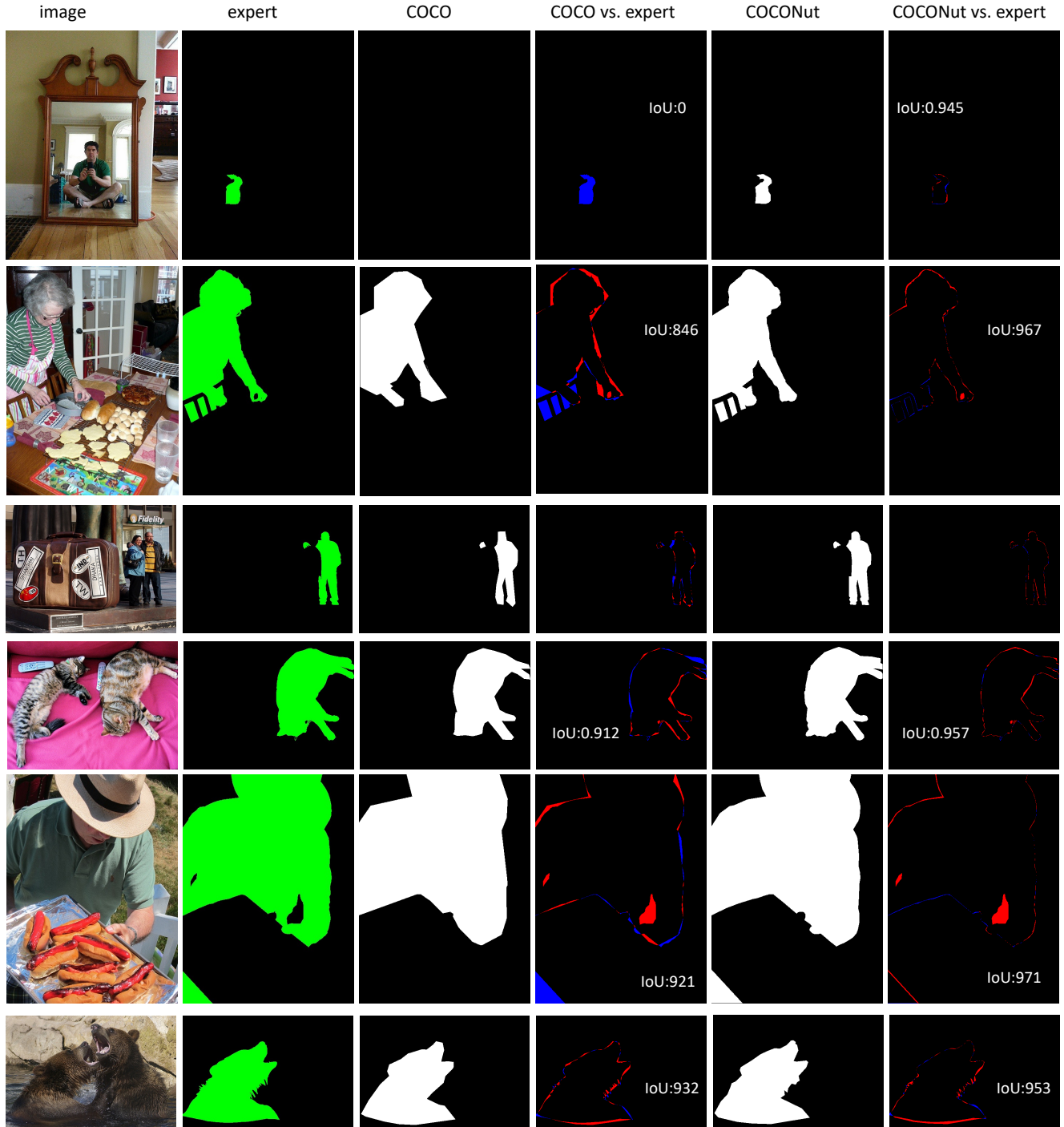
14

Figure 14. **Annotation Comparison:** We show annotations obtained by COCO, COCONut with Box2Mask for 'thing', and our expert rater. COCONut's annotation exhibits sharper boundaries, closely resembling expert results, as evident from higher IoU values. The blue and red regions correspond to extra and missing regions, respectively, compared to the expert mask.

ent from Tab. 6 in the main paper where the mask AP is evaluated using a panoptic segmentation model, we train a

task-specific model on instance segmentation with instance masks only. Tab. 9 summarizes the results. Similar to the

15

| method | backbone | training data | ADE20K-150 | | | A-847 | PC-459 | PC-59 | PAS-21 |
|---|---|---|---|---|---|---|---|---|---|
| | | | PQ | AP$^{mask}$ | mIoU | mIoU | mIoU | mIoU | mIoU |
| FC-CLIP | ConvNeXt-L | COCO | 26.8 | 16.8 | 34.1 | 14.8 | 18.2 | 58.4 | 81.8 |
| | | COCONut-S | 27.3 | 17.3 | 33.8 | 15.3 | 20.4 | 57.5 | 82.1 |
| | | COCONut-B | 27.4 | 17.4 | 33.7 | 15.5 | 20.1 | 58.5 | 82.0 |
| | | COCONut-L | 27.5 | 17.4 | 33.9 | 15.6 | 20.6 | 58.0 | 81.9 |

Table 12. **Benchmarking Open-Vocabulary Segmentation:** We ablate the effect of using different training data to train the mask proposal network of FC-CLIP [69]. The performance is evaluated on multiple segmentation datasets in a zero-shot manner.

| method | backbone | training set | COCO-val AP$^{box}$ | relabeled COCO-val AP$^{box}$ | COCONut-val AP$^{box}$ |
|---|---|---|---|---|---|
| DETA | ResNet50 | COCO | 50.4 | 49.5 | 46.1 |
| | | COCONut-S | 47.8 | 53.8 | 49.5 |
| | | COCONut-B | 50.4 | 54.4 | 51.4 |
| | | COCONut-L | 50.6 | 54.9 | 53.7 |
| | Swin-L | COCO | 59.1 | 58.6 | 56.1 |
| | | COCONut-S | 54.5 | 61.3 | 58.9 |
| | | COCONut-B | 59.3 | 62.2 | 60.1 |
| | | COCONut-L | 60.1 | 62.3 | 61.7 |

Table 13. **Benchmarking Bounding Box Object Detection:** We conduct the experiments using the DETA framework [44], employing various backbones and diverse training and validation sets. The backbones are exclusively pretrained on ImageNet [51].

| method | training set | COCO-val | | relabeled COCO-val | COCONut-val | |
|---|---|---|---|---|---|---|
| | | FID ↓ | mIoU ↑ | FID ↓ mIoU ↑ | FID ↓ | mIoU ↑ |
| GLIGEN | COCO | 18.51 | 32.1 | - 33.7 | 17.4 | 30.9 |
| | COCONut-S | 18.39 | 30.4 | - 34.8 | 16.8 | 32.6 |

Table 14. **Benchmarking Mask-Conditional Image Synthesis:** We conduct the experiments using the GLIGEN framework [35] mIoU is measured with another off-the-shelf Mask2Former [12], as a referee.

findings in panoptic segmentation, we observe consistent improvements across various backbones (ResNet50 [24] and ConvNeXt-L [40]). Additionally, as we increase the size of training dataset, we observe that the improvement gain is decreasing while evaluated on the small COCO-val and relabeled COCO-val set, indicating the performance saturation on the small validation set. By contrast, the proposed COCONut-val set presents a more challenging validation set, where the improvement of stronger backbone and more training images are more noticeable.

**Semantic Segmentation:** We also conduct experiments on training a single semantic segmentation model with semantic masks. Results are shown in Tab. 10. Similar observations are made. We can see subsequent mIoU gains of increasing the training dataset size for semantic specific model. Additionally, we also verify the dataset on semantic segmentation using ViT backbones [14]. We follow the same configuration and use the codebase from ViT-Adapter [9] to conduct our experiments but replace the dataset from COCO-stuff to our COCONut semantic segmentation dataset. As shown in Tab. 11, a similar observation is made: the model saturates when testing on our re-

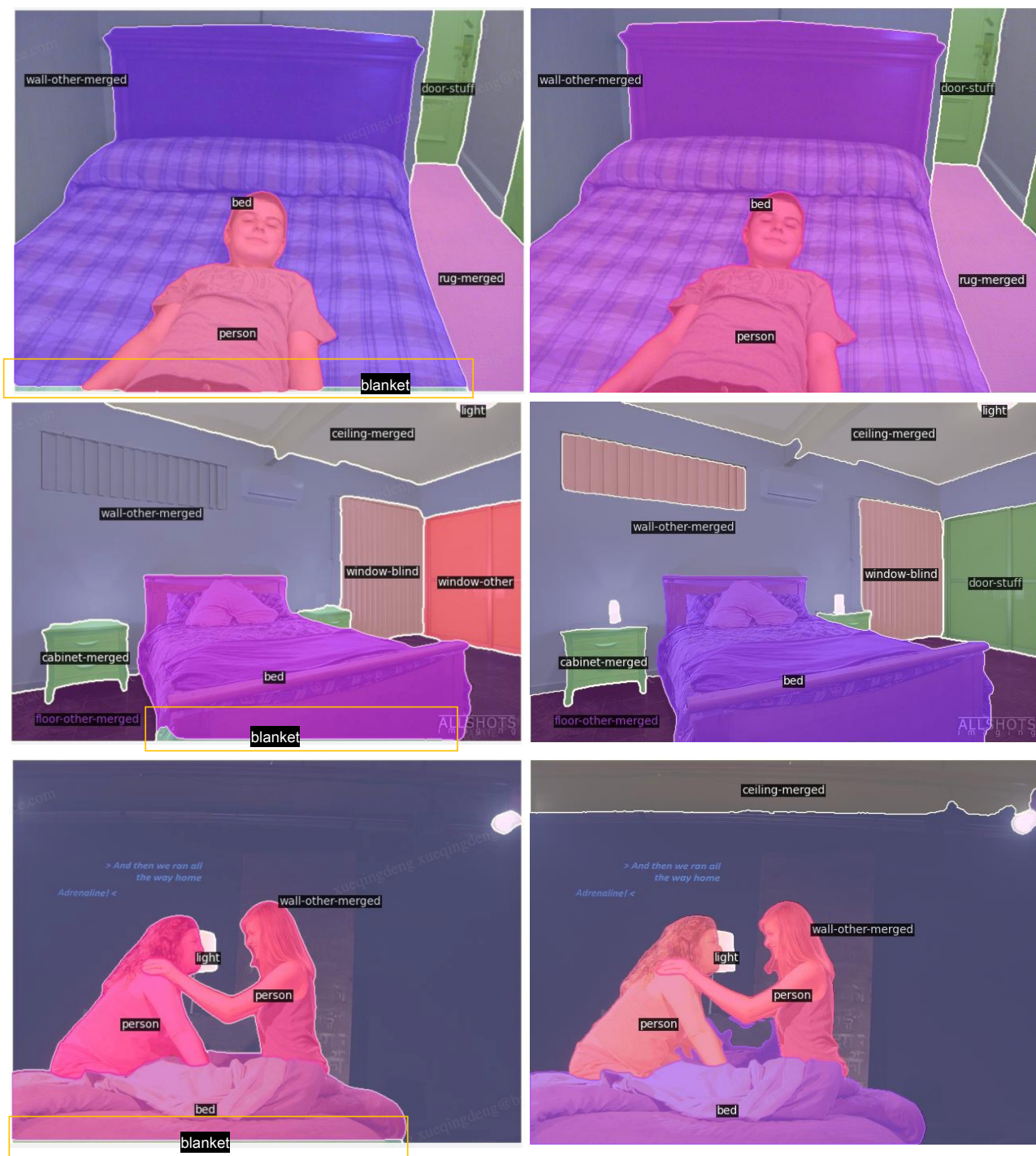labeled COCO-val set but the performance is improved on COCONut-val.

**Open-Vocabulary Segmentation:** Tab. 12 summarizes the results for open-vocabulary segmentation using FC-CLIP. As shown in the table, FC-CLIP benefits from COCONut's high-quality and large-scale annotations, achieving the performance of 27.5 PQ on ADE20K, setting a new state-of-the-art.

**Bounding Box Object Detection:** The results for object detection are shown in Tab. 13. As shown in the table, the detection model with ResNet50 benefits significantly from the high quality data COCONut-S with a large margin of 4.3 on relabeled COCO-val set. Similarly, subsequent gains from training size are observed for both backbones.

**Mask Conditional Image Synthesis:** We conduct mask conditional image synthesis to verify the annotation quality for generation. We employ a mask-conditional model GLIGEN [35] and train the model on paired image-mask data from COCO and COCONut-S separately. Once we have the trained model checkpoint, we perform inference on mask-conditioned generation by giving masks from COCO val set, relabeled COCO-val set, and COCONut-val set individually to compute FID. The lower FID shows better image synthesis performance. Besides, we adopt the off-the-shelf Mask2Former [12] model to perform semantic segmentation by giving the generated images as input and report mIoU for evaluation. As shown in Tab. 14, our high-quality mask annotation can result in better image synthesis with 18.39 FID on COCO-val set and 16.8 FID on COCONut-val set. Besides, the higher-quality generated images can be better inferred via the higher segmentation mIoU scores. Even for a more challenging val set, training on COCONut-S outperforms the COCO dataset.

## C. Label Map Details and Annotation Instruction

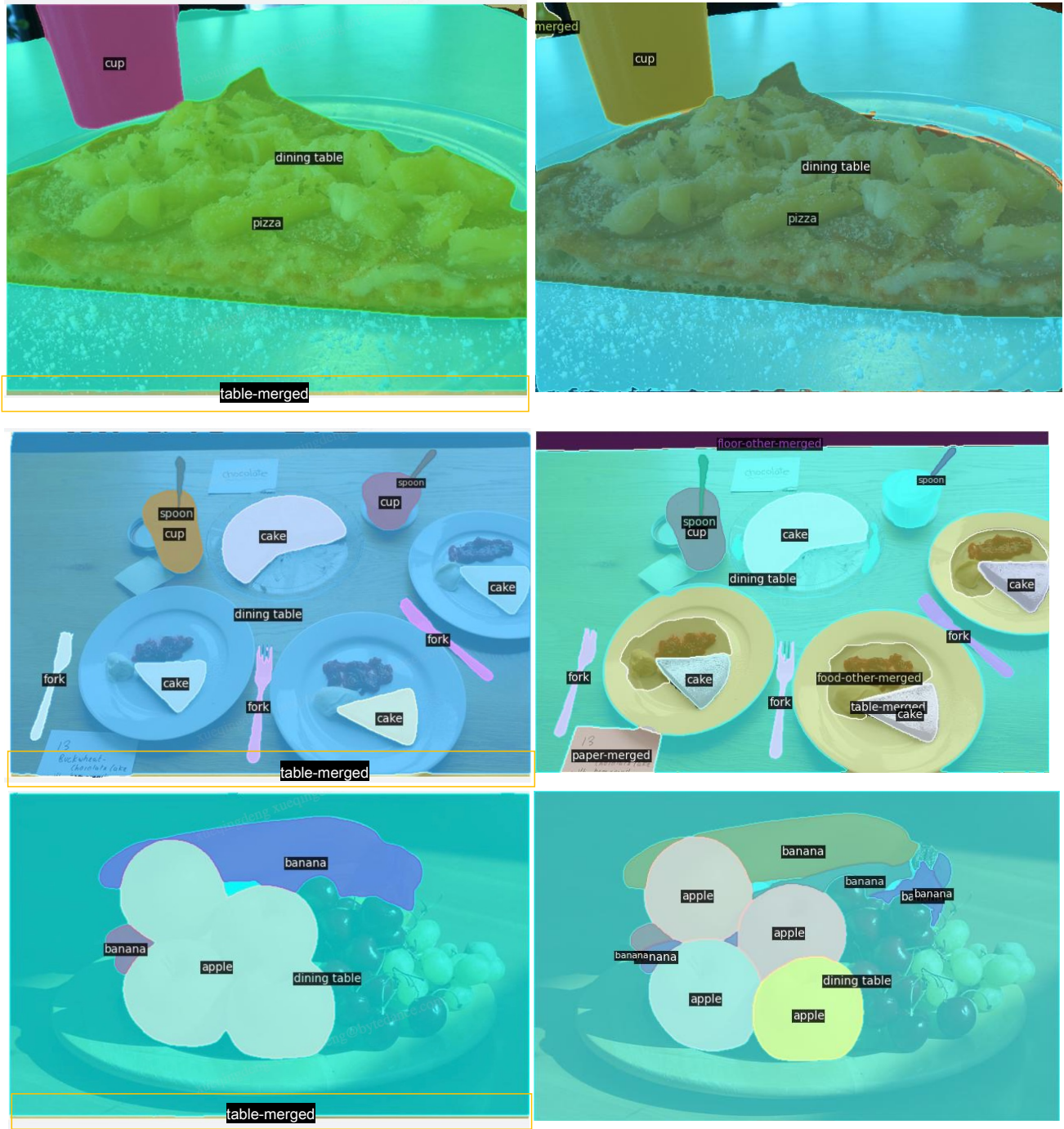Our label map definition strictly follows COCO [36]. However, the COCO definitions of specific categories might be ambiguous. Specifically, we have identified several conflicts between the 'thing' and 'stuff' categories, often resulting in potential mask overlaps. To mitigate the issue, we have meticulously redefined specific categories, detailed in Tab. 15, Tab. 16, and Tab. 17. The definitions of cate-

COCO trained prediction                    COCONut trained prediction

Figure 15. **Influence of Training Data on Predictions:** We present predictions from two models: one trained on original COCO (left) and the other on COCONut (right). The COCO-trained model predicts a small isolated mask, influenced by the biases inherent in the COCO coarse annotations.

COCO trained prediction             COCONut trained prediction

Figure 16. **Influence of Training Data on Predictions:** We present predictions from two models: one trained on original COCO (left) and the other on COCONut (right). The COCO-trained model predicts a small isolated mask, influenced by the biases inherent in the COCO coarse annotations.

gories not included in the tables remain consistent with the original COCO label map.

# References

[1] Radhakrishna Achanta, Appu Shaji, Kevin Smith, Aurelien Lucchi, Pascal Fua, and Sabine Süsstrunk. Slic superpixels

|  | category | COCO definition | COCONut definition |
|---|---|---|---|
| 'thing' | bed | None | A piece of furniture for sleep or rest, typically a framework with a mattress and coverings (from Google dictionary). Thus we will include the pillows, comforter, blanket, and bedding sheets along with the bed frame for bed. |
| 'stuff' | blanket | A loosely woven fabric, used for warmth while sleeping. | As blanket in the bed is included in the category of bed, then we will label blanket on the other surface excluding bed, for example, blanket on the couch or blanket on the bench. |
| 'stuff' | pillow | A rectangular cloth bag stuffed with soft materials to support the head. | To avoid the conflicts from bed in 'thing', we exclude the pillow in the bed while labeling pillow. |
| 'thing' | dining table | None | A table on which meals are served in a dining room (from Google dictionary). In order to have consistent COCO's definition by viewing hundreds of examples, partial table placing with food is also considered as dining table. |
| 'stuff' | table-merged | A piece of furniture with a flat top and one or more legs. | Exclude the cases from dining table aforementioned. Include console table, coffee table, desk and *etc*. |
| 'stuff' | roof | The structure forming the upper covering of a building. | The structure forming the upper covering of a building or vehicle (from Google dictionary). Only the outside coverings will be labeled. COCO also labels the inner side of the coverings while they should be referred as ceiling instead. |
| 'stuff' | house | A smaller size building for human habitation. | A building for human habitation, especially one that is lived in by a family or small group of people (from Google dictionary). Typically it refers to residential house meanwhile residential apartment building is not included. To avoid overlap from roof, a house will not to separated into the parts of roof and the remaining while this happens in COCO. |
| 'stuff' | building-other-merged | Any other type of building or structures. | For the other types of buildings, it consists of diverse types of constructions, for example, churches, stadiums, and *etc*. |
| 'stuff' | wall-tile | A building wall made of tiles, such as used in bathrooms and kitchens. | Follow the same definition from COCO. |
| 'stuff' | wall-stone | A building wall made of stone. | Indoor wall with specific texture of stone and partial outside wall of the building instead of the whole building. In other word, the building built with stone will be labeled as building instead of wall-stone. |
| 'stuff' | wall-wood | A building wall made of wooden material. | Indoor and outside wall made of wood instead of the whole building. |
| 'stuff' | wall-brick | A building wall made of bricks of clay. | Indoor and outside wall made of bricks instead of the whole building. |
| 'stuff' | wall-other-merged | Any other type of wall. | To avoid the conflicts wall categories, we will first label the categories with specific texture and at last we label wall-other-merged. In details, we only label indoor scenes for wall-other-merged, for outdoor scenes, we will use other categories. We also need to exclude the other objects hang on the wall, for example, the frames *etc*. |

Table 15. **Clear Redefinition of Specific COCO Categories:** We present the class definitions by grouping confusing categories for easier comparison to facilitate their distinction (continued in Tab. 16).

| | category | COCO definition | COCONut definition |
|---|---|---|---|
| 'stuff' | gravel | A loose aggregation of small water-worn or pounded stones. | Follow the same definition from COCO. |
| 'stuff' | railroad | A track made of steel rails along which trains run (incl. the wooden beams). | We found that railroad often consists of the gravel and the track. In this scenario, we separate the region of gravel to be labeled as gravel and the remaining parts of the track will be labeled as railroad. |
| 'stuff' | playingfield | A ground marked off for various games (incl. indoor and outdoor). | Follow the same definition. But we found COCO has a large amount of missing masks for playingfield which are mislabeled as dirt-merged instead. We label all the playingfields if they can be identified no matter they are grass based or dirt based grounds. |
| 'stuff' | dirt-merged | Soil or earth (incl. dirt path). | Follow the same definition but exclude dirt-based playingfields. |
| 'stuff' | pavement-merged | A typically raised paved path for pedestrians at the side of a road. | Follow definition from COCO, to be more concrete, it includes side walk. |
| 'stuff' | platform | A paved way leading from one place to another. | COCO does not have consistent labeling masks for platforms while some of them are labeled as pavement-merged. We have a unified definition to take care of these cases. In particular, we label all the paved way for transportation, for example, label the pavement area for the train, subway and *etc.* as platform. |
| 'stuff' | net | An open-meshed fabric twisted, knotted, or woven together at regular intervals. | Follow the same definition but exclude fence made by net. |
| 'stuff' | fence-merged | A thin, human-constructed barrier which separates two pieces of land. | COCO has inconsistent masks for fence-merged and net. We follow a consistent definition to distinguish net from fence when it is not used as a fence to separate two pieces of land. |
| 'thing' | potted plant | None | A plant that is grown in a container, and usually kept inside. There exist masks for flower placed in the vase which contradicting the definition of flower and vase. We exclude these scenarios from potted plant. |
| 'thing' | vase | None | A decorative container, typically made of glass or china and used as an ornament or for displaying cut flowers (google dictionary). |
| 'stuff' | flower | The seed-bearing part of a plant (incl. the entire flower). | COCO does not clarify that whether the flowers that are placed in the vase belong to potted plant or flower. This is confusing when our raters label the images. We give the definition to separate the flower, potted plant and vase. The potted plant will not include any plants which are flowers. Then the potted plant will be labeled together with the plants and pots. While for vase, if the vase has flower, then these parts need to be separate. |

Table 16. **Clear Redefinition of Specific COCO Categories:** We present the class definitions by grouping confusing categories for easier comparison to facilitate their distinction (continued in Tab. 17).

compared to state-of-the-art superpixel methods. *TPAMI*, 34 (11):2274–2282, 2012. 8

[2] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. Coco-stuff: Thing and stuff classes in context. In *CVPR*, 2018. 3,

4, 8

[3] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020. 2

| category | COCO definition | COCONut definition |
|---|---|---|
| food-other-merged | Any other type of food. | To avoid the conflicts from similar categories of 'thing', we explicitly highlight that we DO NOT label those categories. The categories include sandwich (burger), hot dog, pizza, donut, cake, broccoli and carrot. Excluding all the food aforementioned, the other types of food need to be labeled. |
| paper-merged | A material manufactured in thin sheets from the pulp of wood. | Include tissue, toilet paper, poster, kitchen paper towel, and *etc*. They are often shown with a single or multiple pieces of papers. |
| tree-merged | A woody plant, typically having a single trunk growing to a considerable height and bearing lateral branches at some distance from the ground. | Include bush. |
| fruit | The sweet and fleshy product of a tree or other plant. | Exclude fruits in 'thing', banana, orange and apple. Include tomato and all other kinds of fruit. |

Table 17. **Clear Redefinition of Specific COCO Categories:** We clearly redefine certain COCO categories to avoid annotation confusion.

[4] Jieneng Chen, Qihang Yu, Xiaohui Shen, Alan Yuille, and Liang-Chieh Chen. ViTamin: Designing Scalable Vision Models in the Vision-Language Era. *arXiv preprint arXiv:2404.02132*, 2024. 2

[5] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan Yuille. Semantic image segmentation with deep convolutional nets and fully connected crfs. In *ICLR*, 2015. 2

[6] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *TPAMI*, 40(4):834–848, 2017.

[7] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *ECCV*, 2018. 2

[8] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*, 2015. 2, 3

[9] Zhe Chen, Yuchen Duan, Wenhai Wang, Junjun He, Tong Lu, Jifeng Dai, and Yu Qiao. Vision transformer adapter for dense predictions. *arXiv preprint arXiv:2205.08534*, 2022. 16

[10] Bowen Cheng, Maxwell D Collins, Yukun Zhu, Ting Liu, Thomas S Huang, Hartwig Adam, and Liang-Chieh Chen. Panoptic-deeplab: A simple, strong, and fast baseline for bottom-up panoptic segmentation. In *CVPR*, 2020. 2

[11] Bowen Cheng, Ross Girshick, Piotr Dollár, Alexander C. Berg, and Alexander Kirillov. Boundary IoU: Improving object-centric image segmentation evaluation. In *CVPR*, 2021. 8

[12] Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In *CVPR*, 2022. 16

[13] Zheng Ding, Jieke Wang, and Zhuowen Tu. Open-vocabulary panoptic segmentation with maskclip. In *ICML*, 2023. 13

[14] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021. 2, 16

[15] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *IJCV*, 88:303–338, 2010. 3, 13

[16] Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao Nguyen, Ryan Marten, Mitchell Wortsman, Dhruba Ghosh, Jieyu Zhang, et al. Datacomp: In search of the next generation of multimodal datasets. *arXiv preprint arXiv:2304.14108*, 2023. 2

[17] Golnaz Ghiasi, Xiuye Gu, Yin Cui, and Tsung-Yi Lin. Scaling open-vocabulary image segmentation with image-level labels. In *ECCV*, 2022. 2, 13

[18] Ross Girshick. Fast r-cnn. In *ICCV*, 2015. 2

[19] Xiuye Gu, Yin Cui, Jonathan Huang, Abdullah Rashwan, Xuan Yang, Xingyi Zhou, Golnaz Ghiasi, Weicheng Kuo, Huizhong Chen, Liang-Chieh Chen, and David A Ross. Dataseg: Taming a universal multi-dataset multi-task segmentation model. *NeurIPS*, 2023. 2

[20] Rıza Alp Güler, Natalia Neverova, and Iasonas Kokkinos. Densepose: Dense human pose estimation in the wild. In *CVPR*, 2018. 2

[21] Agrim Gupta, Piotr Dollar, and Ross Girshick. LVIS: A dataset for large vocabulary instance segmentation. In *CVPR*, 2019. 3

[22] Bharath Hariharan, Pablo Arbeláez, Ross Girshick, and Jitendra Malik. Simultaneous detection and segmentation. In *ECCV*, 2014. 2, 13

[23] Ju He, Qihang Yu, Inkyu Shin, Xueqing Deng, Alan Yuille,

Xiaohui Shen, and Liang-Chieh Chen. Maxtron: Mask transformer with trajectory attention for video panoptic segmentation. *arXiv preprint arXiv: 2311.18537*, 2023. 2

[24] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 8, 13, 16

[25] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *ICCV*, 2017. 2

[26] Xuming He, Richard S Zemel, and Miguel Á Carreira-Perpiñán. Multiscale conditional random fields for image labeling. In *CVPR*, 2004. 13

[27] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. 9

[28] Lei Ke, Mingqiao Ye, Martin Danelljan, Yifan Liu, Yu-Wing Tai, Chi-Keung Tang, and Fisher Yu. Segment anything in high quality. In *NeurIPS*, 2023. 5

[29] Alexander Kirillov, Ross Girshick, Kaiming He, and Piotr Dollár. Panoptic feature pyramid networks. In *CVPR*, 2019. 2

[30] Alexander Kirillov, Kaiming He, Ross Girshick, Carsten Rother, and Piotr Dollár. Panoptic segmentation. In *CVPR*, 2019. 3, 4, 6, 13

[31] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment anything. In *ICCV*, 2023. 2, 3, 5

[32] Weicheng Kuo, Yin Cui, Xiuye Gu, AJ Piergiovanni, and Anelia Angelova. F-vlm: Open-vocabulary object detection upon frozen vision and language models. In *ICLR*, 2023. 2

[33] 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: Unified image classification, object detection, and visual relationship detection at scale. *IJCV*, 128(7):1956–1981, 2020. 3

[34] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML*, 2023. 2

[35] Yuheng Li, Haotian Liu, Qingyang Wu, Fangzhou Mu, Jianwei Yang, Jianfeng Gao, Chunyuan Li, and Yong Jae Lee. Gligen: Open-set grounded text-to-image generation. In *CVPR*, 2023. 16

[36] 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. 2, 3, 4, 6, 9, 13, 16

[37] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *CVPR*, 2017. 2

[38] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *NeurIPS*, 2023. 2

[39] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, 2021. 6, 13

[40] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *CVPR*, 2022. 2, 6, 8, 13, 16

[41] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015. 2

[42] Qi Lu, Jason Kuen, Shen Tiancheng, Gu Jiuxiang, Guo Weidong, Jia Jiaya, Lin Zhe, and Yang Ming-Hsuan. High-quality entity segmentation. In *ICCV*, 2023. 3, 4

[43] Roozbeh Mottaghi, Xianjie Chen, Xiaobai Liu, Nam-Gyu Cho, Seong-Whan Lee, Sanja Fidler, Raquel Urtasun, and Alan Yuille. The role of context for object detection and semantic segmentation in the wild. In *CVPR*, 2014. 3, 13

[44] Jeffrey Ouyang-Zhang, Jang Hyun Cho, Xingyi Zhou, and Philipp Krähenbühl. Nms strikes back. *arXiv preprint arXiv:2212.06137*, 2022. 3, 4, 6, 13, 16

[45] Dim P Papadopoulos, Jasper RR Uijlings, Frank Keller, and Vittorio Ferrari. Extreme clicking for efficient object annotation. In *ICCV*, 2017. 2

[46] George Papandreou, Tyler Zhu, Liang-Chieh Chen, Spyros Gidaris, Jonathan Tompson, and Kevin Murphy. Personlab: Person pose estimation and instance segmentation with a bottom-up, part-based, geometric embedding model. In *ECCV*, 2018. 2

[47] Siyuan Qiao, Liang-Chieh Chen, and Alan Yuille. Detectors: Detecting objects with recursive feature pyramid and switchable atrous convolution. In *CVPR*, 2021. 2

[48] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021. 2

[49] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *NeurIPS*, 2015. 2

[50] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 2

[51] 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. 16

[52] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *NeurIPS*, 2022. 2

[53] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *NeurIPS*, 2022. 2, 3

[54] 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. 2, 3, 4, 6, 9

[55] Inkyu Shin, Dahun Kim, Qihang Yu, Jun Xie, Hong-Seok Kim, Bradley Green, In So Kweon, Kuk-Jin Yoon, and Liang-Chieh Chen. Video-kmax: A simple unified approach for online and near-online video panoptic segmentation. In *WACV*, 2024. 2

[56] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning for human pose estimation. In *CVPR*, 2019. 2

[57] Shoukun Sun, Min Xian, Fei Xu, Tiankai Yao, and Luca Capriotti. Cfr-icl: Cascade-forward refinement with iterative click loss for interactive image segmentation. *arXiv preprint arXiv:2303.05620*, 2023. 3, 5, 6

[58] Shuyang Sun, Weijun Wang, Andrew Howard, Qihang Yu, Philip Torr, and Liang-Chieh Chen. Remax: Relaxing for better training on efficient panoptic segmentation. *NeurIPS*, 2024. 2

[59] Huiyu Wang, Yukun Zhu, Bradley Green, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. Axial-deeplab: Stand-alone axial-attention for panoptic segmentation. In *ECCV*, 2020.

[60] Huiyu Wang, Yukun Zhu, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. Max-deeplab: End-to-end panoptic segmentation with mask transformers. In *CVPR*, 2021. 2

[61] Weiyun Wang, Min Shi, Qingyun Li, Wenhai Wang, Zhen-hang Huang, Linjie Xing, Zhe Chen, Hao Li, Xizhou Zhu, Zhiguo Cao, et al. The all-seeing project: Towards panoptic visual recognition and understanding of the open world. *arXiv preprint arXiv:2308.01907*, 2023. 3

[62] Mark Weber, Huiyu Wang, Siyuan Qiao, Jun Xie, Maxwell D Collins, Yukun Zhu, Liangzhe Yuan, Dahun Kim, Qihang Yu, Daniel Cremers, Laura Leal-Taixé, Alan Yuille, Florian Schroff, Hartwig Adam, and Liang-Chieh Chen. Deeplab2: A tensorflow library for deep labeling. *arXiv preprint arXiv:2106.09748*, 2021. 13

[63] Jiarui Xu, Sifei Liu, Arash Vahdat, Wonmin Byeon, Xiaolong Wang, and Shalini De Mello. Open-vocabulary panoptic segmentation with text-to-image diffusion models. In *CVPR*, 2023. 2, 13

[64] Xuan Yang, Liangzhe Yuan, Kimberly Wilber, Astuti Sharma, Xiuye Gu, Siyuan Qiao, Stephanie Debats, Huisheng Wang, Hartwig Adam, Mikhail Sirotenko, and Liang-Chieh Chen. Polymax: General dense prediction with mask transformer. In *WACV*, 2024. 2

[65] Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. Coca: Contrastive captioners are image-text foundation models. *TMLR*, 2022. 2

[66] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. *TMLR*, 2022. 2

[67] Qihang Yu, Huiyu Wang, Dahun Kim, Siyuan Qiao, Maxwell Collins, Yukun Zhu, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. Cmt-deeplab: Clustering mask transformers for panoptic segmentation. In *CVPR*, 2022. 2

[68] Qihang Yu, Huiyu Wang, Siyuan Qiao, Maxwell Collins, Yukun Zhu, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. k-means Mask Transformer. In *ECCV*, 2022. 3, 4, 6, 8, 13

[69] Qihang Yu, Ju He, Xueqing Deng, Xiaohui Shen, and Liang-Chieh Chen. Convolutions die hard: Open-vocabulary segmentation with single frozen convolutional clip. *NeurIPS*, 2023. 2, 13, 16

[70] Qihang Yu, Xiaohui Shen, and Liang-Chieh Chen. Towards open-ended visual recognition with large language model. *arXiv preprint arXiv:2311.08400*, 2023. 2

[71] Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun Zhu, Lionel M Ni, and Heung-Yeung Shum. Dino: Detr with improved denoising anchor boxes for end-to-end object detection. In *ICLR*, 2023. 2

[72] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *CVPR*, 2017. 3, 13

[73] Eric Zimmermann, Justin Szeto, Jerome Pasquero, and Frederic Ratle. Benchmarking a benchmark: How reliable is ms-coco? In *ICCV Datacomp Workshop*, 2023. 3, 4