# Instance Segmentation of Scene Sketches Using Natural Image Priors

Mia Tang[1]    Yael Vinker[2]    Chuan Yan[1]    Lvmin Zhang[1]    Maneesh Agrawala[1]

[1]Stanford University                    [2]MIT

{miatang, chuanyan, lvmin, maneesh}@stanford.edu        yaelvink@mit.edu

https://sketchseg.github.io/sketch-seg/

Figure 1. Our method performs instance segmentation of raster sketches. It effectively handles diverse types of sketches, accommodating variations in stroke style and complexity.

## Abstract

*Sketch segmentation involves grouping pixels within a sketch that belong to the same object or instance. It serves as a valuable tool for sketch editing tasks, such as moving, scaling, or removing specific components. While image segmentation models have demonstrated remarkable capabilities in recent years, sketches present unique challenges for these models due to their sparse nature and wide variation in styles. We introduce SketchSeg, a method for instance segmentation of raster scene sketches. Our approach adapts state-of-the-art image segmentation and object detection models to the sketch domain by employing class-agnostic fine-tuning and refining segmentation masks using depth cues. Furthermore, our method organizes sketches into sorted layers, where occluded instances are inpainted, enabling advanced sketch editing applications. As existing datasets in this domain lack variation in sketch styles, we construct a synthetic scene sketch segmentation dataset featuring sketches with diverse brush strokes and varying levels of detail. We use this dataset to demonstrate the robustness of our approach and will release it to promote further research in the field.*

## 1. Introduction

Sketches serve as a powerful tool for visual exploration, ideation, and planning. Traditional sketching workflows often begin with artists working on a single canvas layer (either physical or digital) to maintain creative momentum. As the sketch evolves and requires refinement (e.g., adjustments to composition, perspective, or other elements), artists face the tedious task of manually segmenting different elements of the sketch into discrete, editable layers. Automating the sketch segmentation process offers a promising solution. However, this task presents unique challenges due to the sparse and abstract nature of line drawings, as well as the inherent variability in human sketching styles. Existing methods for scene-level sketch segmentation typically rely on training dedicated models using annotated sketch datasets [5, 26, 47]. However, most available datasets are confined to specific sketch styles and a limited set of object categories, restricting the generalization capabilities of existing methods.

In this work, we introduce SketchSeg, a method for instance segmentation of raster scene sketches that outperforms previous approaches in accommodating a wider variety of sketch styles and concepts. We use the segmenta-

1

tion map to divide the sketch into sorted layers to support effective sketch editing.

Our method builds upon Grounded SAM [24], a state-of-the-art approach for open-vocabulary image segmentation, which has demonstrated remarkable capabilities in segmenting complex scenes across diverse object categories. Grounded SAM combines two models to achieve this: Grounding DINO [17] for object detection and Segment Anything (SAM) [12] for mask generation. We analyze the performance of these models on sketch inputs, revealing that the domain gap between real and sketched objects presents significant challenges for Grounding DINO. In contrast, SAM exhibits a surprising ability to generalize to sketches, though it still faces difficulties specific to the sketch domain. To address the gap in object detection, we fine-tune Grounding DINO on a small subset of annotated scene sketches from the SketchyScene dataset [47]. Our straightforward fine-tuning technique proves highly effective, achieving a substantial improvement in Grounding DINO's detection performance on sketches, with Average Precision increasing from 24% to 75%.

For object segmentation, we apply SAM [12] in the sketch domain, using detected object regions from our fine-tuned Grounding DINO. This is followed by a depth-based refinement stage to resolve ambiguities in overlapping regions. Finally, we decompose the segmented sketch into sorted layers and employ a pretrained image inpainting model [31] to fill in missing regions. This layered representation facilitates sketch editing, allowing users to drag or manipulate segmented objects without the need to manually sketch the affected regions, as we demonstrate in the provided video.

To evaluate our method on diverse scene sketches, we construct a synthetic annotated dataset that extends existing benchmarks along three key dimensions: drawing style, stroke style, and object categories. The dataset integrates two complementary pipelines to enhance diversity. The first pipeline builds on SketchyScene [47], expanding its clipart-like sketches with styles ranging from high-fidelity representations to symbolic, abstract sketches, introducing challenging out-of-distribution cases. The scenes are created in vector format to allow for stroke style variations, including Calligraphic Pen, Charcoal, and Brush Pen styles. The second pipeline leverages the Visual Genome dataset [13], which provides annotated scenes with object variety. Using the InstantStyle method [33], we generate expressive, natural-looking sketches spanning 74 categories, extending SketchyScene's original 45 categories by 54 new categories. Our dataset contains 20,000 annotated scene sketches in total, and is highly extensible. Our evaluations demonstrate that SketchSeg generalizes well to these challenging variations, significantly advancing the state of the art.

## 2. Related Work

### 2.1. Part-Level Sketch Segmentation

The majority of work in the sketch segmentation domain focuses on part-level semantic segmentation, in which the goal is to assign labels to object parts (*e.g.*, the body, wings, and head of a bird). These methods often rely on curated part-level sketch segmentation datasets [6, 8, 10, 14, 36] to train a segmentation model, and use various network architectures, including CNNs [32, 46], RNNs [11, 22, 36], Graph Neural Networks [40, 44], Transformers [34, 45], and more specific techniques such as deformation networks [20] and CRFs [27]. These approaches typically operate on a fixed set of object classes, and recognize a predefined set of object parts within them. Other work focuses on perceptual grouping [14, 15] to achieve class-agnostic segmentation. However all of these methods are designed to tackle part-level segmentation and are not suitable for scene sketches.

### 2.2. Scene-Level Sketch Segmentation

Scene-level sketch segmentation remains largely under-explored. Qi *et al*. [21] extend the perceptual grouping approach to scene-level images, forming semantically meaningful groupings of edges, though with limited accuracy on complicated scenes. Zou *et al*. [47] construct the SketchyScene dataset, providing annotated scene sketches with meaningful layouts of object interactions, and use it to train an instance segmentation model based on the Mask R-CNN architecture [9]. However, their method is limited to the predefined categories included in the dataset, and the proposed dataset contains sketches with clipart-like appearance which challenges the model's ability to generalize to other artistic styles. Building on SketchyScene, Ge *et al*. [7] introduced SKY-Scene and TUB-Scene by replacing its object components with sketches from the Sketchy [26] and TU-Berlin [6] datasets. However, their proposed fusion network is fundamentally limited to the fixed set of classes it was trained on, and the trained network weights are not publicly available. SFSD [43] develops a dataset featuring more complex scene sketches, and utilizes a bidirectional LSTM to produce stroke-level segmentation. Unfortunately, the dataset and model are not publicly available. SketchSeger [38] proposes a hierarchical Transformer-based model for semantic sketch segmentation. However their model is inherently restricted to the predefined set of classes used during training. Bourouis *et al*. [2] finetune the CLIP image encoder [23] on the FS-COCO [5] dataset, leveraging the model's vision-language prior to enable open-vocabulary scene segmentation. However their method is designed for semantic segmentation, and it struggles to generalize to more challenging sketch styles and scene layouts.
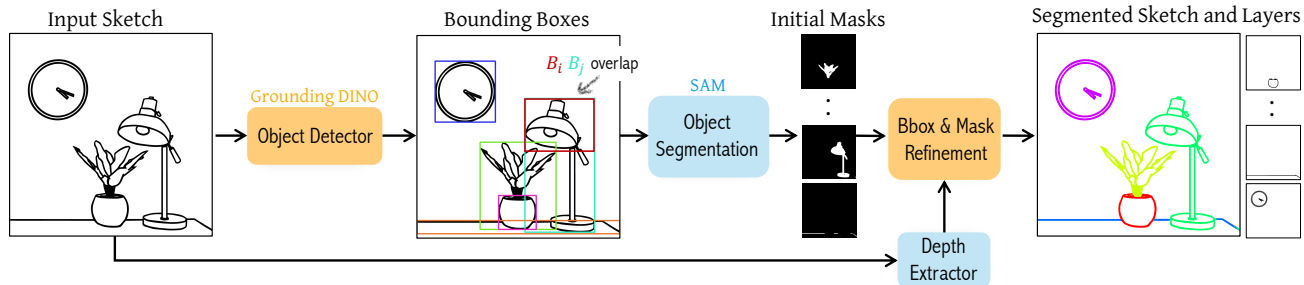
Figure 2. Overview of the sketch segmentation pipeline. Given an input sketch image, our framework first detects bounding boxes using a customized Grounding DINO to obtain region proposals, and then perform segmentation with SAM models. The localization and segmentation are refined by incorporating the depth features. The result segmentation can be viewed as a layered decomposition of object components in the original sketch.

## 2.3. Image Segmentation

The task of image segmentation have been widely explored [1, 4, 9, 35]. The advent of vision-language models [16, 23, 37] has led to numerous object detection and segmentation methods with impressive generalization capabilities [12, 19, 24, 41]. Grounding DINO [17] is a state-of-the-art object detection model trained on over 10 million images. It builds on top of DINO [3], a strong vision encoder, with effective grounding module that fuses visual and textual information, enabling open-vocabulary detection of unseen objects. Segment Anything (SAM) [12] is an image segmentation model trained on over 11 million images and 1.1 billion masks, capable of producing high-quality object masks based on various forms of conditioning such as bounding boxes. Grounded SAM [24], which our method builds upon, combines Grounding DINO and SAM for open-vocabulary image segmentation, achieving robust performance across diverse object categories. Yet, despite demonstrating impressive capabilities on natural images, we show that these models struggle with segmenting sketches.

## 3. Method

Given a raster sketch, our goal is to produce a segmentation map such that pixels belonging to the same object instance are grouped together. Based on the segmentation map, we also divide the sketch into layers, sorted by depth. Our pipeline is illustrated in Figure 2. Given the input sketch, we first perform object detection using a fine-tuned Grounding DINO model, which produces a set of candidate object bounding boxes. These bounding boxes are then used to produce an initial set of object masks with a pre-trained Segment Anything (SAM) [12] model. Next, we perform a refinement stage that leverages scene depth information and classical morphological operations to assign the final segmentation. This stage also employs a pre-trained inpainting model [31] to produce scene layers.

### 3.1. Sketch-Aware Object Detection

Grounding DINO [17] is an object detection model which outputs bounding boxes for recognized object instances, based on a given text prompt describing the scene. While effective for natural images, the model in its original configuration demonstrates limited generalization to sketches (we show this numerically in Section 5). To address this limitation, we fine-tune Grounding DINO on sketches. The largest available annotated sketch dataset containing complex scenes is SketchyScene [47]. It contains 30K segmented sketches across 45 class labels. We find that a naive fine-tuning with the SketchyScene data leads to severe overfitting to the small set of predefined object classes.

To overcome this overfitting, we propose a class-agnostic fine-tuning strategy. Instead of relying on predefined class labels, we train the model to distinguish between instances based on their visual characteristics, aiming to push the model to rely on Gestalt properties such as closure, continuity, and emergence, to group together strokes forming a single object. Specifically, we utilize a small subset of 5000 sketches from the SketchyScene dataset, and consolidate their class labels into a single label, "object". We use a Grounding DINO model initialized with a pretrained Swin Transformer [18] backbone and fine-tune the model's detection head for bounding box prediction. For training, we employ standard object detection losses used in the original Grounding DINO training (Focal Loss, L1 Loss, and GIoU Loss), while eliminating the class recognition loss. At inference, the model is prompted with the input image and the word "object" to detect all potential object instances in the scene. This results in an initial set of $k$ bounding boxes $B = \{B_i\}_{i=1}^k$ and a confidence score per bounding box.

### 3.2. Mask Extraction and Bounding Boxes Refinement

Once the bounding boxes are obtained, we use a pretrained Segment Anything (SAM) model [12] to extract masks for
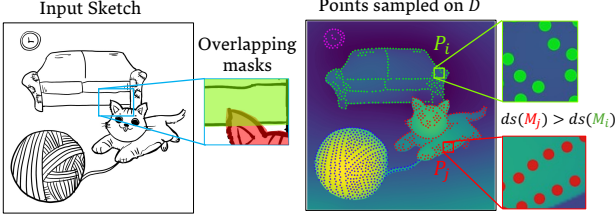
Figure 3. Resolving ambiguities in overlapping regions. The depth map $D$ is sampled along sketch pixels at evenly spaced points, and the sampled points are grouped by their corresponding object (e.g., $P_i$ corresponds to the $i$'th object). Each object is assigned a depth score based on the majority of depth values from the sampled points. Ambiguous pixels are then assigned to the mask with the highest depth score, prioritizing foreground objects.
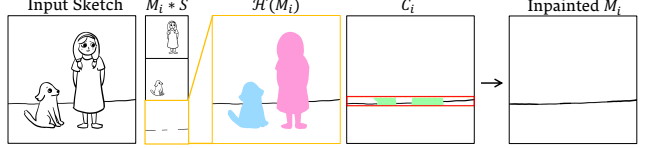


Figure 4. Object layers are isolated and inpainted using a pre-trained SDXL model. The inpainting mask for each object is defined by intersecting overlapping masks with the object's bounding box.

the corresponding objects directly from the sketch. This results in an initial set of $k$ masks $M = \{M_i\}_{i=1}^{k}$ which we refine using simple binary operations such as morphological closing and flood-fill, to eliminate small artifacts.

Next, we use the refined masks to enhance the set of generated bounding boxes. A common practice is to eliminate redundant bounding boxes often corresponding to the same object (such as $B_i, B_j$ shown in red and blue in Figure 2) using Non-Maximum Suppression (NMS), which filters out bounding boxes with low confidence scores that has significant intersection with others. However, IoU of the bounding boxes may not reliably reflect object overlap in cases where objects do not fully cover the pixels in their bounding boxes. This issue is especially pronounced in sketches, which are sparser than photorealistic images. We use the initial set of masks to compute a more fine-grained IoU. Specifically, for a pair of overlapping bounding boxes $B_i$ and $B_j$, we extract the regions within the bounding boxes that intersect with the sketch: $M_i * S, M_j * S$, and compute the IoU of these regions to define an "overlapping" score between two objects $i, j$:

$$\mathcal{O}(i,j) = \text{IoU}(M_i * S, M_j * S). \qquad (1)$$

For an overlapping pair $B_i, B_j$ if $\mathcal{O}(i,j) > 0.5$, we consider the detections to be covering the same object and retain only the bounding box with the highest confidence score. This results in a filtered set of bounding boxes $\hat{B} \subseteq B$ and their corresponding masks $\hat{M} \subseteq M$.

The filtered set of masks may still include overlapping regions, as illustrated in Fig. 3, where it is unclear which instance the pixels should be associated with. To resolve such overlaps and assign each pixel to a single object instance, we give priority to objects in the foreground. We utilize DepthAnything [39] to extract the depth map $D$ of the input sketch. We then sample $D$ along the sketch pixels at equally spaced points $P = \{p_1, p_2, \ldots, p_n\}$, analogous to projecting rays through the sketch pixels to the scene. For each mask $M_i$, we identify the subset of points that lie

within the mask: $P_i = \{p \in P | p \in M_i\}$. For example, in Figure 3, $P_i$ represents the set of points belonging to the sofa, while $P_j$ denotes the set of points belonging to the cat. For each point $p \in P_i$ we associate a depth value $D(p)$ using the depth map. We then compute a depth score for each mask as the mode of the depth values associated with its sampled points:

$$ds(M_i) = \arg\max_{D(p)} \text{count}(\{D(p)|p \in P_i\}). \qquad (2)$$

Based on this score, we assign ambiguous pixels to the mask with the highest depth score, ensuring that foreground objects take precedence. Lastly, to ensure complete coverage of the sketch, we employ a watershed-based [28] refinement, propagating existing mask labels to previously unlabeled sketch pixels.

### 3.3. Layer Inpainting

As a final step we extract complete layers for each object in the sketch, inpainting any occluded regions using a pre-trained SDXL inpainting model [31]. The goal of this stage is to support basic sketch editing operations, such as translation and scaling. We isolate each object $i$ by intersecting the sketch with its corresponding mask $M_i * S$ (Figure 13). We then identify the group of masks that intersect with $M_i$: $\mathcal{H}(M_i) = \{M_j | M_j \cap M_i \neq \varnothing\}$. Finally, we define the inpainting mask $C_i$ as the intersection of $\mathcal{H}(M_i)$ with the object's bounding box: $C_i = \mathcal{H}(M_i) \cap B_i$ (shown in green in Fig. 13), and feed it into the pretrained inpainting model.

### 3.4. Implementation Details

Optimization is performed with the AdamW optimizer, configured with an initial learning rate of 6e-5 and a weight decay of 0.0005 to promote generalization and prevent overfitting. Training is conducted with a batch size of 4 and automatic learning rate scaling to ensure stable updates and efficient adaptation. For the train, validation, and test sets, we sampled 5,000, 500, and 500 images, respectively, from the original SketchyScene train, val, and test splits. The experiment was conducted on a single NVIDIA 4090 GPU, with a total training time of four hours.
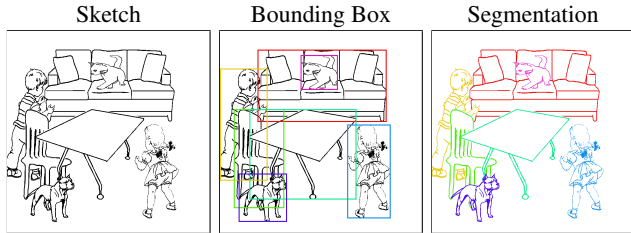
| Sketch | Bounding Box | Segmentation |
|---|---|---|

Figure 5. SketchyScene dataset provides ground truth object bounding boxes and pixel-level instance segmentation masks for scene layouts.

## 4. Scene Sketch Segmentation Benchmark

To evaluate our performance across a diverse set of sketches, we construct a synthetic annotated scene-sketch dataset. This dataset focuses on three key axes of variation, designed to extend existing datasets: (1) drawing style, (2) stroke style, and (3) object categories. We define drawing style as a spectrum ranging from symbolic, which emphasizes abstraction and simplified representation, to realistic, which prioritizes detailed and lifelike depiction. We define stroke variation as the differences in texture, width, and flow that characterize individual strokes, similar to the variety of brush types in digital drawing software. Our dataset combines two complementary pipelines to enhance diversity and object variety.

**SketchyScene Layouts** The SketchyScene dataset [47] consists of 7,265 scene layouts containing 45 object categories. These layouts are of high quality, as they were manually constructed by humans. Each data sample includes an input sketch, object class labels, bounding boxes, and a pixel-wise segmentation map (see Fig. 5). The sketches in the dataset share a consistent clipart-like style. We extend the SketchyScene dataset to include more diverse sketch styles and stroke variations. Specifically, we incorporate recent object sketching methods that introduce significantly different sketch appearances compared to SketchyScene. These include CLIPasso [29], which transforms images of individual objects into sketches with relatively high image fidelity, and SketchAgent [30], which generates symbolic sketches resembling children's drawings, offering a more challenging out-of-distribution case. Both techniques produce vector-format sketches, which we use to assemble scene sketches while avoiding artifacts caused by transformations. Each object is placed at its ground truth location and scaled to fit its bounding box while preserving its aspect ratio. Figure 6 demonstrate sketches produced from a given SketchyScene layout. Exploring stroke variation is crucial for testing the robustness of automatic segmentation approaches, as real-world scenarios often involve highly diverse sketch styles. We augment the vector sketch us-

ing three distinct brush styles through the Adobe Illustrator Scripting API - Calligraphic Pen, Charcoal, and Brush Pen. For each brush type, we manually select the stroke width that best preserved a natural and visually appealing result.

**Extended Categories** To extend the range of 45 object categories available in SketchyScene, we utilize the Visual Genome dataset [13], a large-scale dataset containing diverse and richly annotated images containing over 33,877 distinct object categories. We use InstantStyle [33], a state-of-the-art style transfer method, to generate corresponding raster sketches from the input scene images, and segment the sketch objects based on the provided image segmentation. As sketches are typically sparse, and very small objects may disappear during the translation from image to sketch, we filtered the dataset to include 1068 images containing five to ten distinct objects per scene. Our dataset expands the class categories of SketchyScene by introducing 54 additional object classes, containing in total 72 categories. A few examples of the resulting dataset are shown in Figure 7.

## 5. Results

Figures 1, 8, 11 and 24 present qualitative results of our method across a diverse range of sketches. These include various object categories, both abstract and detailed scenes, different styles, and sketches from our new dataset featuring stroke variations and challenging abstractions. Our method effectively handles object categories beyond those used in our fine-tuning from the SketchyScene dataset, such as toys, furniture, and food items. Our approach successfully addresses challenging scenarios, such as detailed scenes with numerous objects and occluded objects, as seen in Figure 1 and the first row of Figure 8. More results are provided in the supplementary material.

### 5.1. Comparisons

We evaluate our method alongside existing scene sketch segmentation approaches, including SketchyScene [47] and the method proposed by Bourouis et al. [2]. Additionally, we include Grounding DINO [17] as a baseline, applying it directly to sketches, using Recognize-Anything Model (RAM) [42] for automatic labeling. Our evaluation dataset consists of 7746 samples: 1,113 test samples from the SketchyScene dataset, 1,113 samples from each of the CLIPasso brush styles: Calligraphic Pen, Charcoal, and Brush Pen, 1,113 samples from the SketchAgent style, and 1068 samples from the InstantStyle sketch dataset, spanning 99 object categories. Each method was applied to these datasets following their recommended best practices. Note that the SketchyScene mask generation implementation relies on legacy dependencies that are no longer executable.

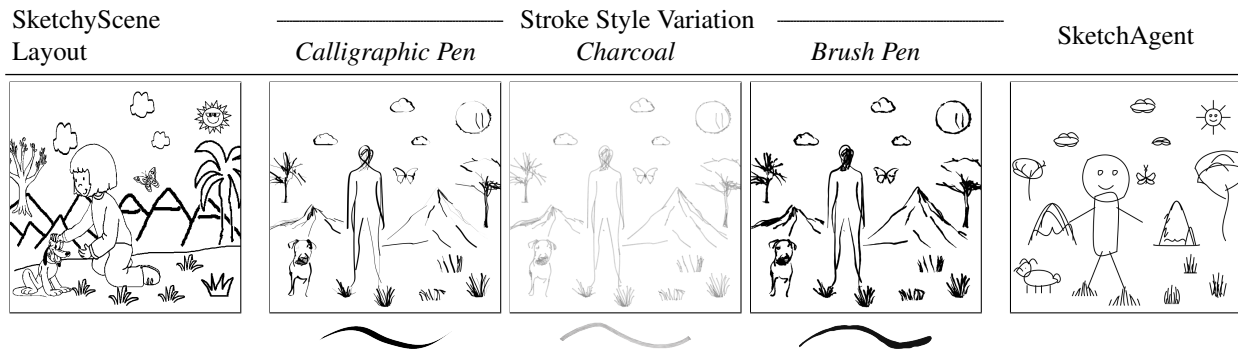| SketchyScene Layout | ——————— Stroke Style Variation ——————— | | | SketchAgent |
| | *Calligraphic Pen* | *Charcoal* | *Brush Pen* | |

Figure 6. Samples from our synthetic dataset. We augment the SketchyScene dataset by generating vector sketches with varied drawing styles based on SketchyScene's scene layouts. Stroke style variation is introduced by re-rendering the scenes with three different brush styles. Additionally, we create a more symbolic and challenging sketch type, resembling children's drawings, shown on the right, while maintaining the same scene layouts.

Table 1. Quantitative comparisons for object detection. We report IoU, AR, and AP metrics across seven datasets, along with the mean and standard deviation for each method. IoU measures the overlap between predicted and ground-truth bounding boxes, AR evaluates the ability to detect all relevant objects, and AP combines precision and recall across varying IoU thresholds from 50% to 95%. AP@50 and AP@75 indicate Average Precision at IoU thresholds of 50% and 75%, respectively, reflecting stricter requirements for bounding box overlap. Our method demonstrates consistent improvements across nearly all datasets and metrics, significantly outperforming baselines, especially in detecting sketch objects with precision.

| | IoU ↑ | | | AR ↑ | | | AP ↑ | | | AP@50 ↑ | | | AP@75 ↑ | | |
| | SketchyS | G-DINO | **Ours** | SketchyS | G-DINO | **Ours** | SketchyS | G-DINO | **Ours** | SketchyS | G-DINO | **Ours** | SketchyS | G-DINO | **Ours** |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SketchyScene | 0.55 | 0.27 | **0.72** | 0.42 | 0.27 | **0.86** | 0.36 | 0.24 | **0.83** | 0.79 | 0.31 | **0.93** | 0.17 | 0.27 | **0.88** |
| SketchAgent | 0.27 | 0.16 | **0.73** | 0.19 | 0.16 | **0.79** | 0.15 | 0.12 | **0.75** | 0.39 | 0.18 | **0.87** | 0.05 | 0.13 | **0.78** |
| C-Base | 0.48 | 0.31 | **0.80** | 0.35 | 0.30 | **0.85** | 0.29 | 0.26 | **0.83** | 0.70 | 0.37 | **0.93** | 0.11 | 0.30 | **0.87** |
| C-Calligraphic | 0.48 | 0.28 | **0.72** | 0.35 | 0.28 | **0.84** | 0.29 | 0.24 | **0.81** | 0.71 | 0.34 | **0.93** | 0.11 | 0.27 | **0.86** |
| C-Charcoal | 0.40 | 0.27 | **0.76** | 0.29 | 0.27 | **0.84** | 0.24 | 0.24 | **0.79** | 0.59 | 0.33 | **0.94** | 0.08 | 0.27 | **0.88** |
| C-BrushPen | 0.41 | 0.27 | **0.75** | 0.30 | 0.27 | **0.82** | 0.25 | 0.23 | **0.79** | 0.60 | 0.33 | **0.90** | 0.10 | 0.26 | **0.82** |
| InstantStyle | 0.20 | **0.49** | 0.45 | 0.17 | 0.48 | **0.61** | 0.12 | 0.37 | **0.51** | 0.29 | 0.53 | **0.69** | 0.08 | 0.40 | **0.52** |
| All | 0.39 | 0.26 | **0.70** | 0.29 | 0.29 | **0.80** | 0.24 | 0.24 | **0.75** | 0.58 | 0.34 | **0.88** | 0.11 | 0.27 | **0.80** |
| | ±0.12 | ±0.04 | ±0.11 | ±0.08 | ±0.08 | ±0.08 | ±0.08 | ±0.07 | ±0.11 | ±0.18 | ±0.10 | ±0.08 | ±0.03 | ±0.07 | ±0.12 |



Input Image    InstantStyle Sketch    Instance Segmentation

Figure 7. Illustration of our synthetic dataset. The input images are sourced from the Visual Genome dataset [13], which we filter to a subset of scenes containing 5 to 10 object instances. We generate the corresponding sketches with InstantStyle [33].

Table 2. Quantitative comparisons from image segmentation. We report Accuracy and IoU metrics across seven datasets, along with the mean and standard deviation for each method. Our method consistently outperforms baselines across all datasets.

| Metric | Acc ↑ | | | IoU ↑ | | |
| Dataset | SketchyS | G-SAM | **Ours** | SketchyS | G-SAM | **Ours** |
|---|---|---|---|---|---|---|
| SketchyScene | 0.79 | 0.53 | **0.92** | 0.72 | 0.26 | **0.88** |
| SketchAgent | 0.39 | 0.35 | **0.88** | 0.38 | 0.16 | **0.84** |
| C-Base | 0.70 | 0.54 | **0.91** | 0.64 | 0.32 | **0.88** |
| C-Calligraphic | 0.66 | 0.50 | **0.87** | 0.63 | 0.30 | **0.86** |
| C-Charcoal | 0.59 | 0.47 | **0.85** | 0.43 | 0.26 | **0.84** |
| C-BrushPen | 0.66 | 0.51 | **0.89** | 0.54 | 0.29 | **0.85** |
| InstantStyle | 0.43 | 0.65 | **0.70** | 0.32 | 0.44 | **0.78** |
| All | 0.60 | 0.50 | **0.86** | 0.52 | 0.29 | **0.78** |
| | ±0.14 | ±0.09 | ±0.07 | ±0.14 | ±0.08 | ±0.03 |

Therefore, we used SAM for mask generation based on their detected bounding boxes, which yielded better performance
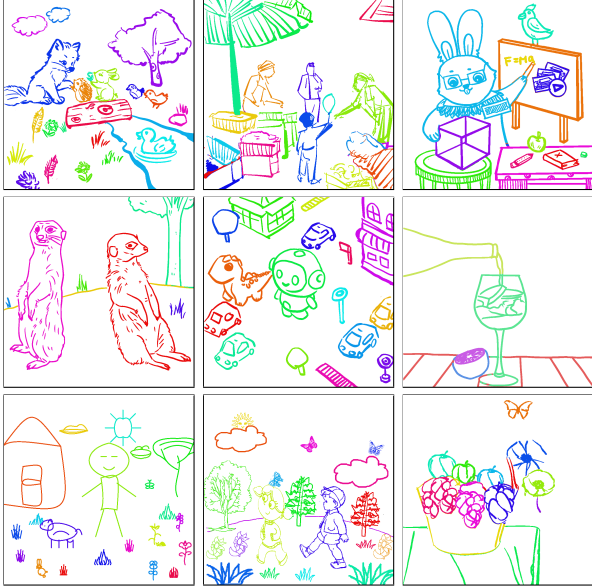
Figure 8. Sketch segmentation results obtained by our method for a diverse set of sketch styles and levels of complexity.

Table 3. Quantitative comparisons from image segmentation on the filtered datasets. We report Accuracy and IoU metrics across seven datasets, along with the mean and standard deviation for each method. OpenVocab performs semantic segmentation, and requires an input text prompt. We provide ground truth class labels as input prompts to generate segmentations.

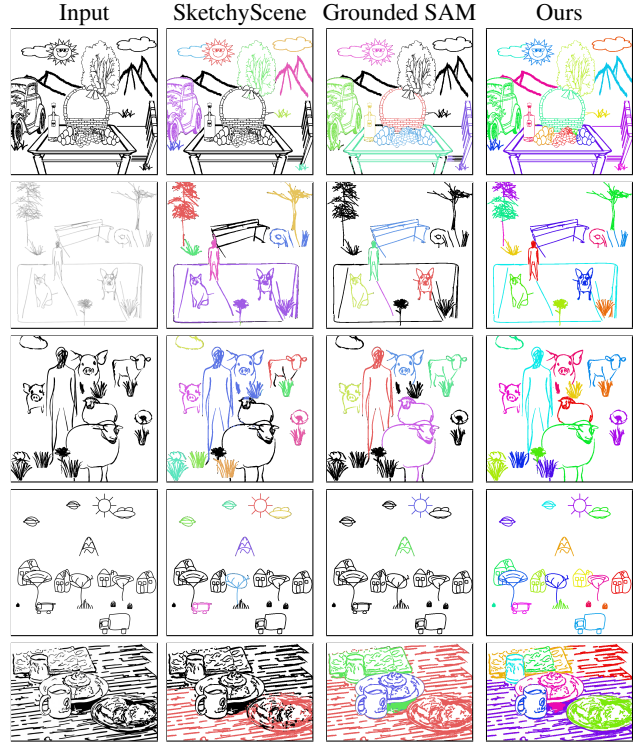| Metric Dataset | Acc ↑ | | IoU ↑ | |
|---|---|---|---|---|
| | OpenVocab | **Ours** | OpenVocab | **Ours** |
| SketchyScene | 0.25 | **0.93** | 0.16 | **0.90** |
| SketchAgent | 0.23 | **0.93** | 0.16 | **0.88** |
| C-Base | 0.29 | **0.94** | 0.19 | **0.92** |
| C-Calligraphic | 0.32 | **0.88** | 0.22 | **0.88** |
| C-Charcoal | 0.55 | **0.85** | 0.43 | **0.85** |
| C-BrushPen | 0.26 | **0.91** | 0.12 | **0.88** |
| InstantStyle | 0.26 | **0.79** | 0.16 | **0.73** |
| All | 0.30 | **0.89** | 0.20 | **0.86** |
| | ±0.11 | ±0.05 | ±0.10 | ±0.06 |



Figure 9. Qualitative comparison of instance segmentation methods. Each row corresponds to a different sketch style, with the first row illustrating a sample from SketchyScene and the remaining rows showcasing samples from our dataset. Black pixels indicate regions where segmentation was not applied. Our method effectively segments sketch pixels into distinct instances without introducing artifacts, outperforming alternative approaches.

than reported in the original paper. Figure 9 illustrates selected segmentation results for all instance segmentation methods, with additional results provided in the supplementary material. Since the method by Bourouis *et al.* [2] is designed for semantic segmentation rather than instance segmentation, it was evaluated separately on a filtered subset of our dataset, where each scene contained only one instance per class.

**Object Detection Evaluation** For object detection, we report Intersection over Union (IoU), which measures the overlap between predicted and ground-truth bounding boxes, Average Recall (AR), which evaluates the ability to detect all relevant objects, and Average Precision (AP), which combines precision and recall across various IoU thresholds. Our goal is to assess the ability to precisely detect any object in the sketch, regardless of its class. To achieve this, we calculate the mean of these metrics across object instances rather than across classes. The results for instance segmentation methods are summarized in Tab. 1. The SketchyScene method performs well on the SketchyScene data, while its performance significantly declines across all metrics when applied to other datasets, particularly for the challenging styles of the SketchAgent samples. In contrast, our method demonstrates consistent performance across all datasets, achieving an average AR score of 0.8, as shown in the last row of the table. Notably, our method outperforms SketchyScene even on its native dataset, demonstrating the effectiveness of leveraging priors from pretrained models on natural images and adapting them for sketch segmentation. The scores obtained for our baseline method, Grounding DINO, support our claim that this model struggles to generalize to the domain of sketches

7

without adaptation, despite its strong performance on natural images. This is evident from the large margin in scores between our method and Grounding DINO, seeing an increase of 44% in IoU, 51% increase in AR, and 51% increase in AP. Furthermore, while we fine-tuned Grounding DINO exclusively on the SketchyScene dataset, the results in the table confirm that this approach surprisingly generalizes well to very different types of sketches and object categories.

**Segmentation Evaluation**    We evaluate the final segmentation results using two common metrics: Pixel Accuracy (Acc), which measures the ratio of correctly labeled pixels to the total pixel count in a sketch, and Intersection over Union (IoU), which evaluates the overlap between the predicted and ground-truth segmentation masks. The results for instance segmentation methods are presented in Tab. 2, while the results for Bourouis *et al.* [2] are shown separately in Tab. 3 since it performs semantic segmentation and requires dataset filtering. As shown, our method outperforms alternative approaches across both metrics, with a particularly notable advantage over Grounded SAM. Additionally, our method demonstrates robustness to various styles, especially excelling on the challenging SketchAgent style compared to other methods.
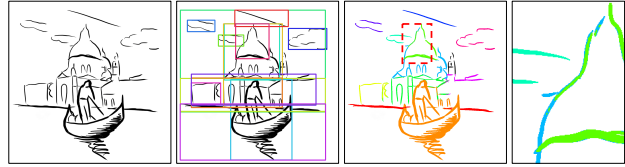
# 6. Limitations and Future Work

While our method successfully segments scene sketches across various styles and challenging cases, it has some limitations. First, our bounding box filtering technique may still include undesired boxes, potentially introducing artifacts when combining masks into the final segmentation (Figure 10a). Second, our mask generation relies on SAM, which generally produces good masks but can occasionally introduce artifacts, particularly for objects occupying large regions in the sketch (Figure 10b). Even after applying our refinement stage, some artifacts may persist in the final segmentation. Future work could address this issue by fine-tuning SAM specifically for sketches or incorporating a learned refinement stage.

# 7. Conclusions

We introduced SketchSeg, a method for instance segmentation of raster scene sketches. Our approach adapts Grounding DINO, an object detection model trained on natural images, to the sketch domain through class-agnostic fine-tuning. We utilized Segment Anything (SAM) for segmentation along with a refinement stage that incorporates depth cues to resolve ambiguous pixels. Our method significantly improves upon state-of-the-art approaches in this domain, demonstrating the utility of natural image priors for sketch understanding tasks. We additionally provide a synthetic

(a) Too many bounding boxes



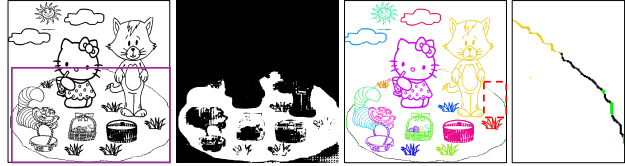(b) Artifacts in masks, especially for large bounding boxes



Figure 10. Examples showcasing the limitations of our approach. (a) The bounding box filtering process can still retain undesired boxes, leading to artifacts in the final segmentation. (b) SAM masks oftentimes are decent, but sometimes do contain noticeable artifacts for larger objects.

scene-level annotated sketch dataset encompassing a wide range of object categories and significant variations in drawing styles. Our experiments demonstrate that SketchSeg is robust to these variations, achieving consistent performance across diverse datasets.

# References

[1] Daniel Bolya, Chong Zhou, Fanyi Xiao, and Yong Jae Lee. Yolact: Real-time instance segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019. 3

[2] Ahmed Bourouis, Judith E Fan, and Yulia Gryaditskaya. Open vocabulary semantic scene sketch understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4176–4186, 2024. 2, 5, 7, 8

[3] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *ICCV*, 2021. 3

[4] Bowen Cheng, Ishan Misra, Alexander G. Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. 2022. 3

[5] Pinaki Nath Chowdhury, Aneeshan Sain, Ayan Kumar Bhunia, Tao Xiang, Yulia Gryaditskaya, and Yi-Zhe Song. Fscoco: Towards understanding of freehand sketches of common objects in context. In *ECCV*, 2022. 1, 2

[6] Mathias Eitz, James Hays, and Marc Alexa. How do humans sketch objects? *ACM Trans. Graph. (Proc. SIGGRAPH)*, 31(4):44:1–44:10, 2012. 2

[7] Ce Ge, Haifeng Sun, Yi-Zhe Song, Zhanyu Ma, and Jianxin Liao. Exploring local detail perception for scene sketch semantic segmentation. *IEEE Transactions on Image Processing*, 31:1447–1461, 2022. 2

[8] Songwei Ge, Vedanuj Goswami, C. Lawrence Zitnick, and Devi Parikh. Creative sketch generation, 2020. 2

[9] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn, 2018. 2, 3

[10] Zhe Huang, Hongbo Fu, and Rynson W. H. Lau. Data-driven segmentation and labeling of freehand sketches. *ACM Trans. Graph.*, 33(6), Nov. 2014. 2

[11] Kurmanbek Kaiyrbekov and Metin Sezgin. Deep stroke-based sketched symbol reconstruction and segmentation. *IEEE Computer Graphics and Applications*, 40(1):112–126, 2020. 2

[12] 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. *arXiv:2304.02643*, 2023. 2, 3

[13] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Li Fei-Fei. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *Int. J. Comput. Vision*, 123(1):32–73, May 2017. 2, 5, 6, 1, 3

[14] Ke Li, Kaiyue Pang, Jifei Song, Yi-Zhe Song, Tao Xiang, Timothy M. Hospedales, and Honggang Zhang. Universal perceptual grouping, 2018. 2

[15] Ke Li, Kaiyue Pang, Yi-Zhe Song, Tao Xiang, Timothy M. Hospedales, and Honggang Zhang. Toward deep universal sketch perceptual grouper. *IEEE Transactions on Image Processing*, 28:3219–3231, 2019. 2

[16] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*, 2023. 3

[17] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023. 2, 3, 5

[18] 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. 3

[19] Matthias Minderer, Alexey Gritsenko, Austin Stone, Maxim Neumann, Dirk Weissenborn, Alexey Dosovitskiy, Aravindh Mahendran, Anurag Arnab, Mostafa Dehghani, Zhuoran Shen, Xiao Wang, Xiaohua Zhai, Thomas Kipf, and Neil Houlsby. Simple open-vocabulary object detection with vision transformers, 2022. 3

[20] Anran Qi, Yulia Gryaditskaya, Tao Xiang, and Yi-Zhe Song. One sketch for all: One-shot personalized sketch segmentation. *CoRR*, abs/2112.10838, 2021. 2

[21] Yonggang Qi, Yi-Zhe Song, Tao Xiang, Honggang Zhang, Timothy Hospedales, Yi Li, and Jun Guo. Making better use of edges via perceptual grouping. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1856–1865, 2015. 2

[22] Yonggang Qi and Zheng-Hua Tan. Sketchsegnet+: An end-to-end learning of rnn for multi-class sketch semantic segmentation. *IEEE Access*, 7:102717–102726, 2019. 2

[23] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. 2, 3

[24] Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kunchang Li, He Cao, Jiayu Chen, Xinyu Huang, Yukang Chen, Feng Yan, Zhaoyang Zeng, Hao Zhang, Feng Li, Jie Yang, Hongyang Li, Qing Jiang, and Lei Zhang. Grounded sam: Assembling open-world models for diverse visual tasks, 2024. 2, 3

[25] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2021. 1

[26] Patsorn Sangkloy, Nathan Burnell, Cusuh Ham, and James Hays. The sketchy database: learning to retrieve badly drawn bunnies. *ACM Trans. Graph.*, 35(4), July 2016. 1, 2

[27] Rosália G. Schneider and Tinne Tuytelaars. Example-based sketch segmentation and labeling using crfs. *ACM Trans. Graph.*, 35(5), July 2016. 2

[28] L. Vincent and P. Soille. Watersheds in digital spaces: an efficient algorithm based on immersion simulations. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 13(6):583–598, 1991. 4

[29] Yael Vinker, Ehsan Pajouheshgar, Jessica Y. Bo, Roman Christian Bachmann, Amit Haim Bermano, Daniel Cohen-Or, Amir Zamir, and Ariel Shamir. Clipasso: Semantically-aware object sketching. *ACM Trans. Graph.*, 41(4), jul 2022. 5, 1

[30] Yael Vinker, Tamar Rott Shaham, Kristine Zheng, Alex Zhao, Judith E Fan, and Antonio Torralba. Sketchagent: Language-driven sequential sketch generation, 2024. 5, 1

[31] Patrick von Platen, Suraj Patil, Anton Lozhkov, Pedro Cuenca, Nathan Lambert, Kashif Rasul, Mishig Davaadorj, Dhruv Nair, Sayak Paul, William Berman, Yiyi Xu, Steven Liu, and Thomas Wolf. Diffusers: State-of-the-art diffusion models. https://github.com/huggingface/diffusers, 2022. 2, 3, 4

[32] Fei Wang, Shujin Lin, Hanhui Li, Hefeng Wu, Tie Cai, Xiaonan Luo, and Ruomei Wang. Multi-column point-cnn for sketch segmentation. *Neurocomputing*, 392:50–59, 2020. 2

[33] Haofan Wang, Matteo Spinelli, Qixun Wang, Xu Bai, Zekui Qin, and Anthony Chen. Instantstyle: Free lunch towards style-preserving in text-to-image generation. *ArXiv*, abs/2404.02733, 2024. 2, 5, 6, 1

[34] Jiawei Wang and Changjian Li. Contextseg: Sketch semantic segmentation by querying the context with attention. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3679–3688, 2024. 2

[35] Xinlong Wang, Rufeng Zhang, Chunhua Shen, Tao Kong, and Lei Li. SOLO: A simple framework for instance segmentation. *IEEE T. Pattern Analysis and Machine Intelligence (TPAMI)*, 2021. 3

[36] Xingyuan Wu, Yonggang Qi, Jun Liu, and Jie Yang. Sketchsegnet: A rnn model for labeling sketch strokes. *2018 IEEE 28th International Workshop on Machine Learning for Signal Processing (MLSP)*, pages 1–6, 2018. 2

[37] Bin Xiao, Haiping Wu, Weijian Xu, Xiyang Dai, Houdong Hu, Yumao Lu, Michael Zeng, Ce Liu, and Lu Yuan. Florence-2: Advancing a unified representation for a variety of vision tasks. *arXiv preprint arXiv:2311.06242*, 2023. 3

[38] Jie Yang, Aihua Ke, Yaoxiang Yu, and Bo Cai. Scene

sketch semantic segmentation with hierarchical transformer. *Knowledge-Based Systems*, 280:110962, 2023. 2

[39] Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything: Unleashing the power of large-scale unlabeled data. In *CVPR*, 2024. 4

[40] Lumin Yang, Jiajie Zhuang, Hongbo Fu, Xiangzhi Wei, Kun Zhou, and Youyi Zheng. Sketchgnn: Semantic sketch segmentation with graph neural networks. *ACM Trans. Graph.*, 40(3), Aug. 2021. 2

[41] Haotian* Zhang, Pengchuan* Zhang, Xiaowei Hu, Yen-Chun Chen, Liunian Harold Li, Xiyang Dai, Lijuan Wang, Lu Yuan, Jenq-Neng Hwang, and Jianfeng Gao. Glipv2: Unifying localization and vision-language understanding. *arXiv preprint arXiv:2206.05836*, 2022. 3

[42] Youcai Zhang, Xinyu Huang, Jinyu Ma, Zhaoyang Li, Zhaochuan Luo, Yanchun Xie, Yuzhuo Qin, Tong Luo, Yaqian Li, Shilong Liu, et al. Recognize anything: A strong image tagging model. *arXiv preprint arXiv:2306.03514*, 2023. 5

[43] Zhengming Zhang, Xiaoming Deng, Jinyao Li, Yukun Lai, Cuixia Ma, Yongjin Liu, and Hongan Wang. Stroke-based semantic segmentation for scene-level free-hand sketches. *Vis. Comput.*, 39(12):6309–6321, Dec. 2022. 2

[44] Yixiao Zheng, Jiyang Xie, Aneeshan Sain, Zhanyu Ma, Yi-Zhe Song, and Jun Guo. Ende-gnn: An encoder-decoder gnn framework for sketch semantic segmentation. In *2022 IEEE International Conference on Visual Communications and Image Processing (VCIP)*, pages 1–5, 2022. 2

[45] Yixiao Zheng, Jiyang Xie, Aneeshan Sain, Yi-Zhe Song, and Zhanyu Ma. Sketch-segformer: Transformer-based segmentation for figurative and creative sketches. *IEEE Transactions on Image Processing*, 32:4595–4609, 2023. 2

[46] Xianyi Zhu, Yi Xiao, and Yan Zheng. Part-level sketch segmentation and labeling using dual-cnn. In *International Conference on Neural Information Processing*, 2018. 2

[47] Changqing Zou, Qian Yu, Ruofei Du, Haoran Mo, Yi-Zhe Song, Tao Xiang, Chengying Gao, Baoquan Chen, and Hao Zhang. Sketchyscene: Richly-annotated scene sketches. In *ECCV*, pages 438–454. Springer International Publishing, 2018. 1, 2, 3, 5
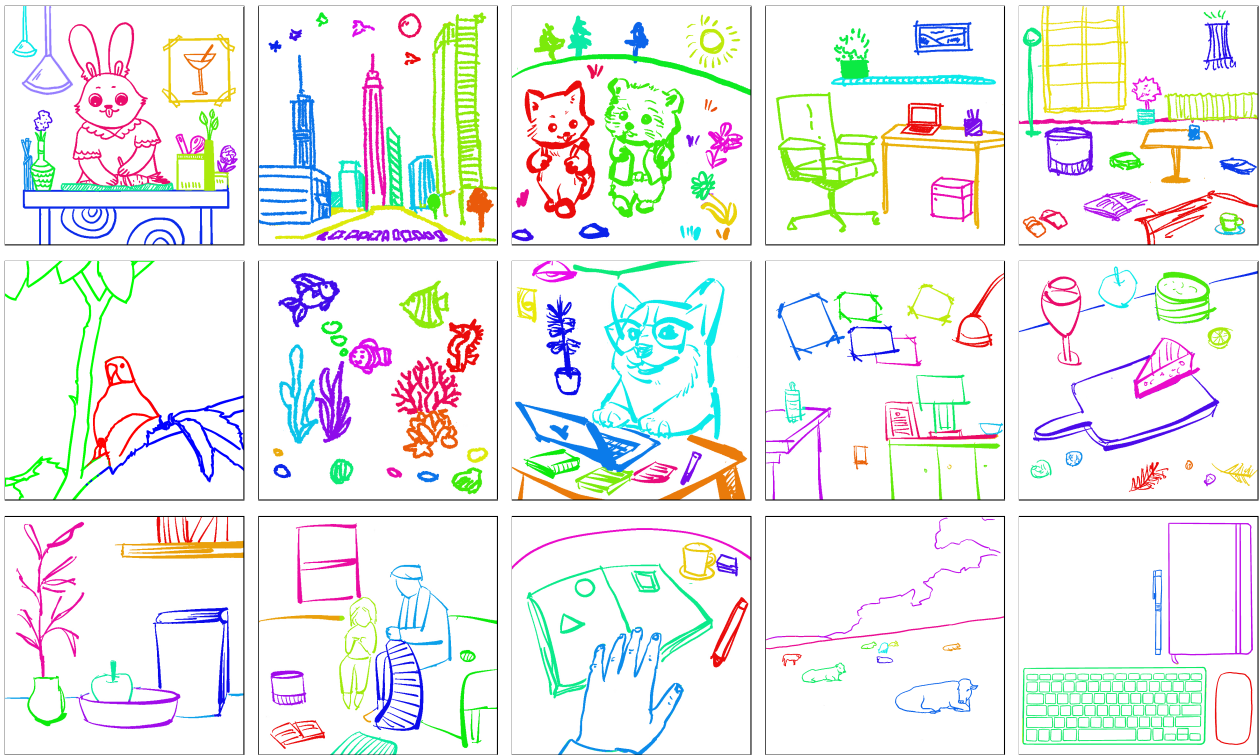
Figure 11. Scene sketch segmentations produced by our method for diverse types of scenes and sketch styles.

Figure 12. SketchSeg segmentation results on a diverse set of sketches from the SketchyScene dataset (first row) and our proposed dataset, which includes three types of stroke style variations (second, third, and fourth rows), challenging abstract sketches generated with SketchAgent (fifth row), and detailed sketches generated with InstantStyle (last row).

# Instance Segmentation of Scene Sketches Using Natural Image Priors
## Supplementary Material



Figure 13. Interactive interface for sketch editing, powered by our instance segmentation and layer completion algorithm.

# Table of Contents

## A. Sketch Editing Interface

Our sketch segmentation and layering technique facilitates sketch editing, allowing users to drag or manipulate segmented objects without the need to manually sketch the affected regions. We demonstrate this through an interactive sketch editing interface (Figure 13) that enables users to upload a sketch, which is then segmented and transformed into completed, ordered layers as detailed in our paper. This facilitates more efficient sketch editing by allowing artists to easily move, copy, or delete pixels associated with specific object instances, as the sketch is represented as an ordered list of layers. **Please see demo video for more examples.**

## B. Synthetic Dataset

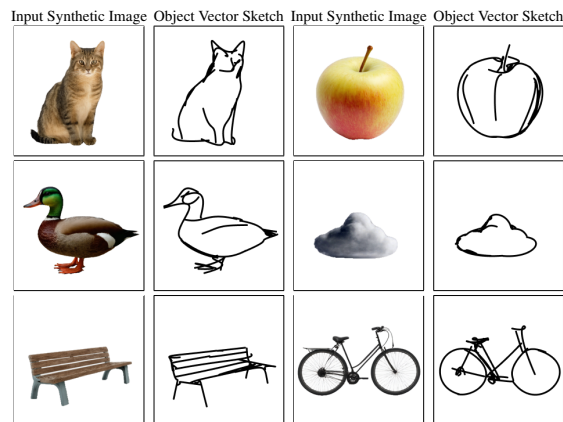In this section, we provide additional details on our synthetic data creation process.



Figure 14. Examples pairs of input synthetic image and output generated object vector sketch.

### B.1. Generating Vector Scene Sketches

We employ CLIPasso [29] and SketchAgent [30] to generate diverse vector sketches of single objects. For each generation method, we create 10 distinct object instances for all 45 classes in the SketchyScene dataset [47], ensuring sufficient variability in our synthetic scenes. For CLIPasso's image-to-vector conversion, we first generate photorealistic synthetic images using SDXL [25]. The generation process uses a consistent prompt template: *"A realistic image of a {class_name} with a blank background"*. Figure 14 demonstrates representative pairs of synthetic input images and their corresponding generated vector sketches. For SketchAgent, we generate sketches directly from class labels as text prompts, producing 10 samples per class. Figure 15 illustrates representative examples of the generated object sketches.

### B.2. Generating Sketches from Natural Images

To expand beyond SketchyScene's object categories, we employ InstantStyle [33] to transform a subset of Visual Genome [13] images into sketches, using a single CLIPasso object sketch as the style reference. Figure 16 showcases a gallery of examples from our InstantStyle-generated scene sketch dataset.

## C. Additional Qualitative Comparisons

We present additional qualitative comparisons between our approach and baseline methods across benchmark scene sketch datasets in Fig. 17 for SketchyScene, Fig. 18 for SketchAgent, Fig. 19 for CLIPasso, Fig. 20 for InstantStyle, to accompany our numerical evaluations included in the paper. To compare with semantic segmentation method,
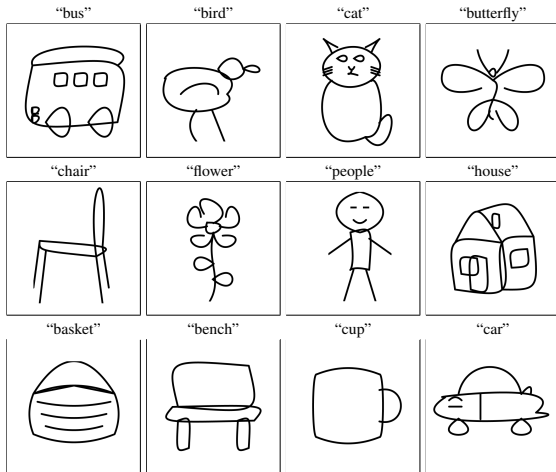
Figure 15. Examples of object vector sketches generated by SketchAgent.

namely OpenVocab, by Bourouis *et al*. [2], we created a filtered version of all seven datasets, where each scene contains at most one instance per object class. These filtered scenes remain challenging for existing methods despite their reduced complexity. We show qualitative results in Fig. 21 for filtered SketchyScene dataset, Fig. 22 for filtered SketchAgent dataset, and Fig 23 for filtered CLI-Passo dataset, to accompany our qualitative results shown in the paper.

Figure 16. **Example sketches from our InstantStyle dataset**. These sketches are derived from Visual Genome [13] containing 5 to 10 annotated objects. For visual clarity, we mask unsegmented regions in the generated sketches. This dataset contains 54 new categories beyond SketchyScene's original 45 classes, and contain 1068 sketches in total.
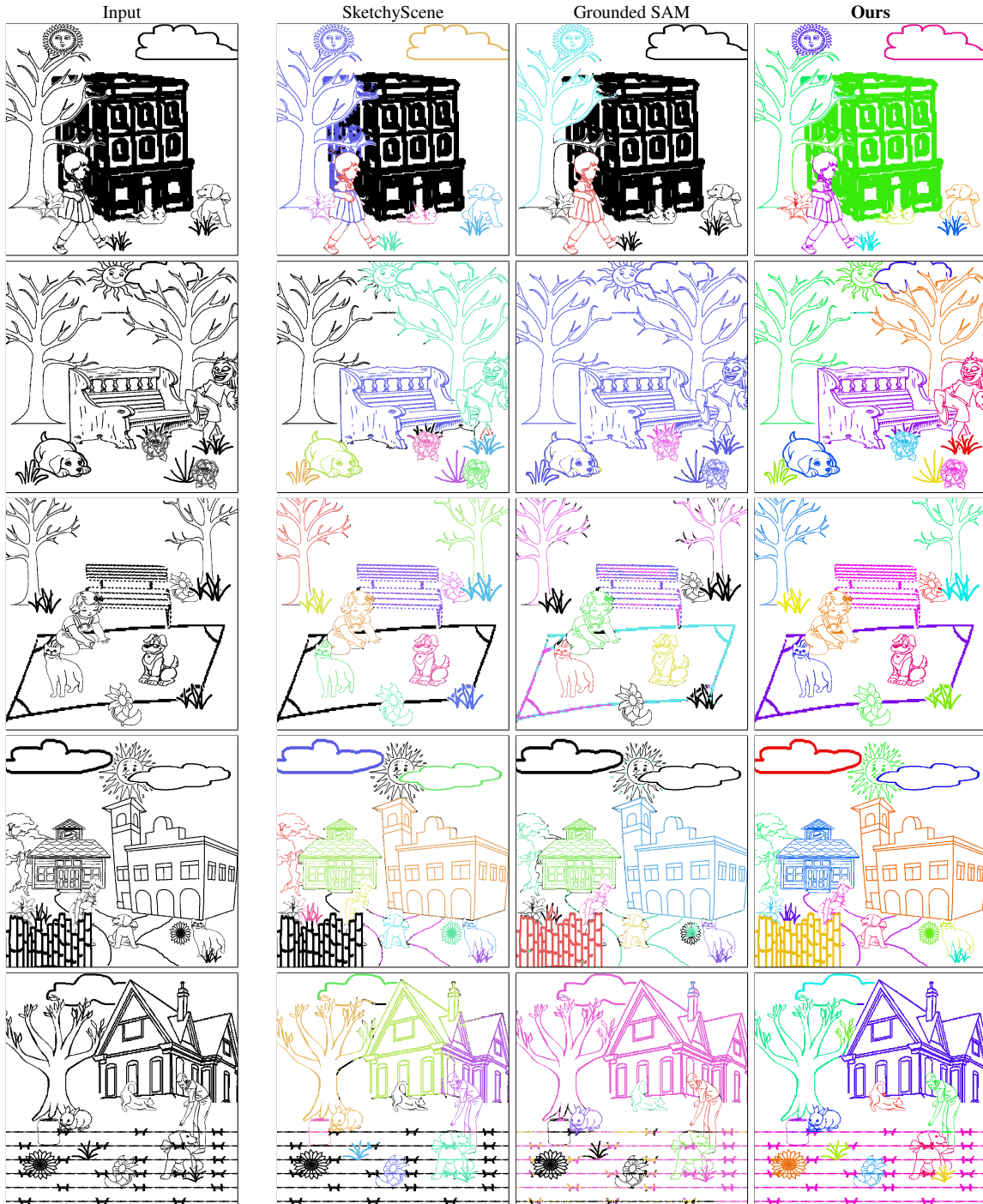
Figure 17. **Qualitative comparison of instance segmentation methods on the SketchyScene dataset.** Our method surpasses SketchyScene and Grounded SAM by delivering fine-grained, semantically consistent segmentations with precise boundaries. In the urban scene (row 1), our approach accurately segments all object instances, cleanly separating the sun, tree, building, and character from each other. For the park scenes (rows 2 and 3), it captures intricate details, such as the dog's face and picnic items, which are either missed or over-segmented by other methods. In the residential and cottage scenes (rows 4 and 5), our method effectively delineates repetitive patterns like fences and handles dense objects like trees, preserving structural integrity where other approaches struggle. These results highlight the robustness of our method in managing complex and detailed sketches.

Figure 18. **Qualitative comparison of instance segmentation methods on the SketchAgent dataset.** Both SketchyScene and Grounded SAM struggle to segment these abstract sketches, which differ significantly from their training data of clipart-like and real objects. Our method successfully segments individual instances while maintaining object boundaries, even in challenging cases like the last row where objects overlap with a grid-patterned picnic blanket.
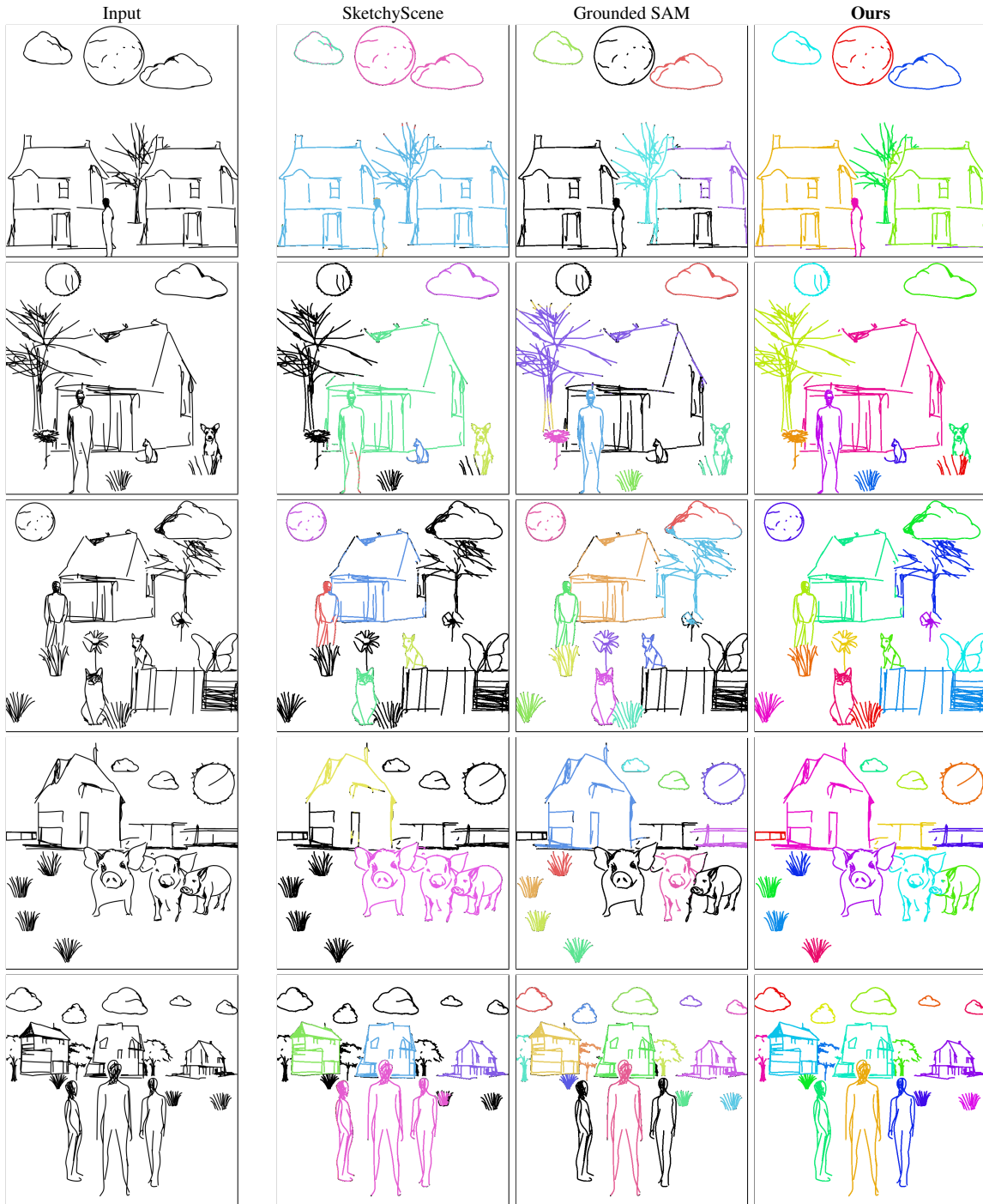
Figure 19. **Qualitative comparison of instance segmentation methods on the CLIPasso base dataset.** Our method successfully detects both large and small objects with ambiguous openings in their silhouettes, whereas baseline methods either merge multiple instances into one or fail to detect them entirely.

6

Figure 20. **Qualitative comparison of instance segmentation methods on the InstantStyle dataset.** This dataset poses a significantly greater challenge compared to the others due to its increased complexity, including diverse perspectives, intricate textures, and frequent occlusions. Despite these difficulties, our method effectively locates object instances, such as the glass of water in row 2, the fork in row 3, and the food in dishes and bottles in row 4. Even with ambiguous shapes, as shown in row 5, our method outperforms GroundedSAM by successfully segmenting the umbrella on the right separately from the person holding it. In the final row, our approach demonstrates its capability by accurately segmenting both the pile of clothes on the couch and the couch itself.
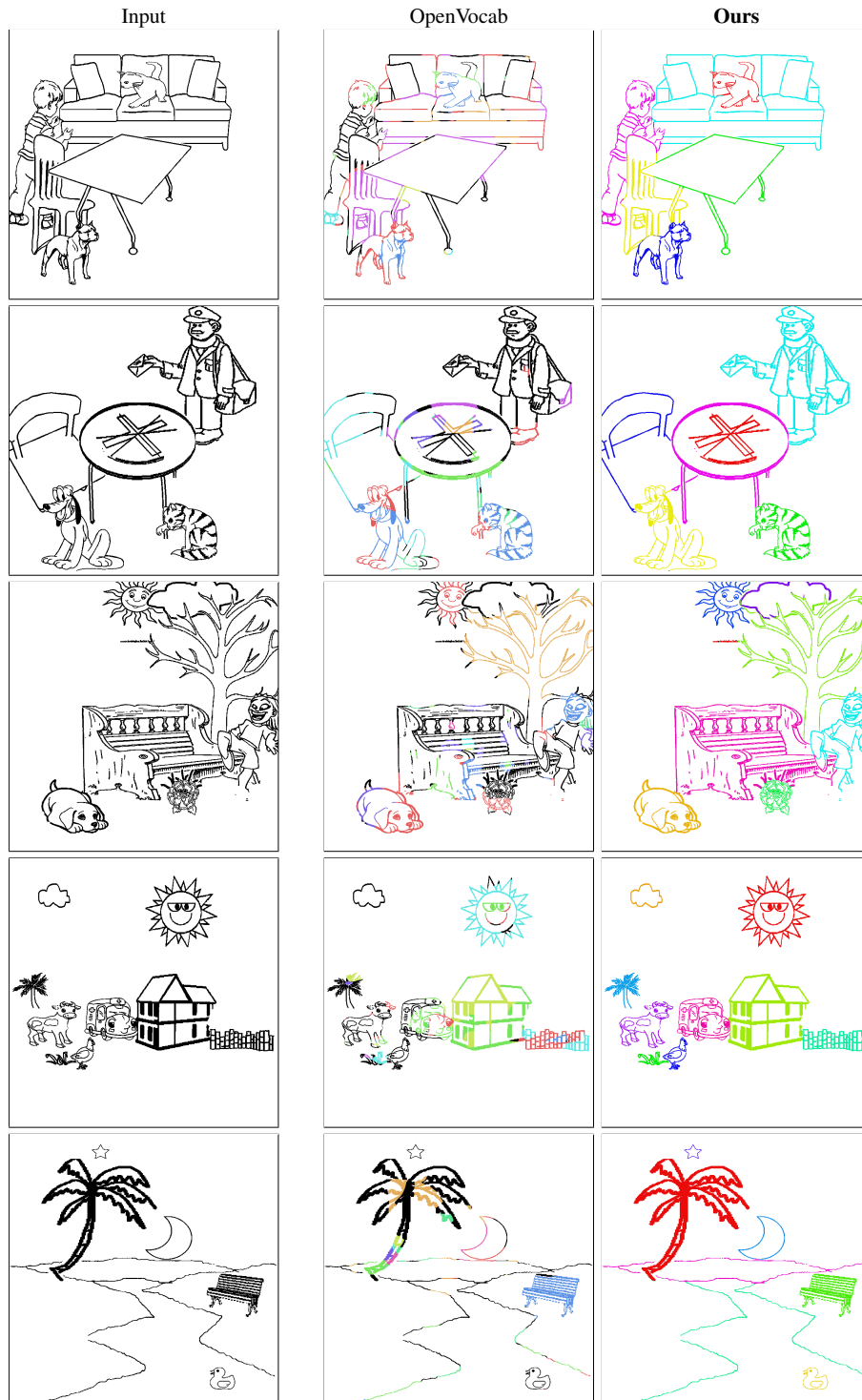
Figure 21. **Qualitative comparison of segmentation on filtered SketchyScene dataset.** We prompt the semantic segmentation model with ground truth class labels; however, it struggles to locate the objects due to the significant deviation of this sketch style from the data it was trained on. Our performance on both complete and filtered scenes are equally as robust, detecting and segmenting object instances precisely.
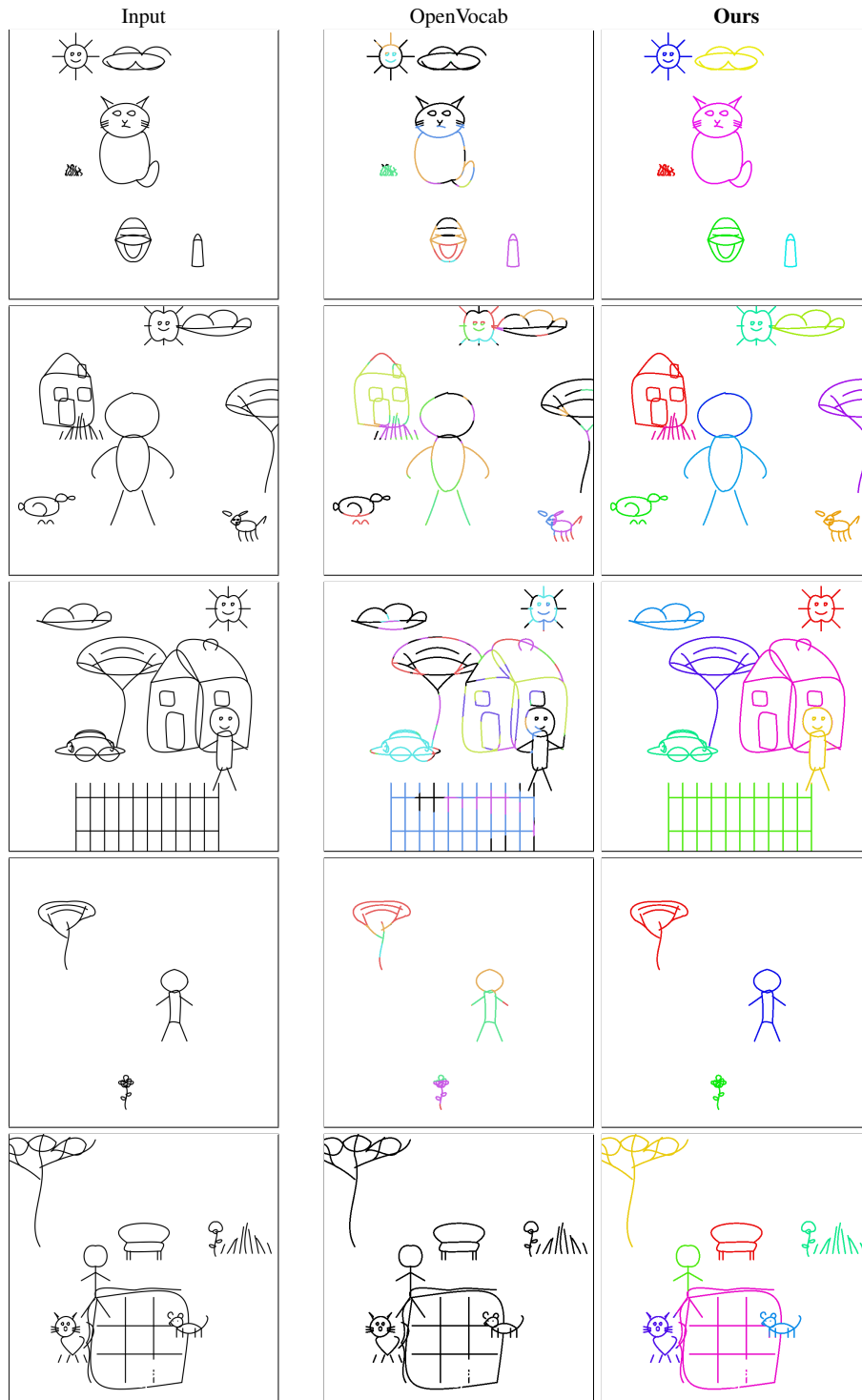
8

Figure 22. **Qualitative comparison of segmentation on filtered SketchAgent dataset.** We prompt the semantic segmentation model with ground truth class labels; however, OpenVocab struggles to accurately identify instances of the correct class, often assigning multiple class labels to the same object, as indicated by the color gradients. Sometimes it is unable to label any sketch pixels in the scene (row 5). In contrast, our method effectively segments object instances, ensuring clear separation and consistent labeling.

9

Figure 23. **Qualitative comparison of segmentation on filtered CLIPasso dataset.** We prompt the semantic segmentation model with ground truth class labels; however, OpenVocab struggles to accurately identify instances of the correct class, often assigning multiple class labels to the same object, as indicated by the color gradients. Sometimes it is unable to label any sketch pixels in the scene (row 4). In contrast, our method effectively segments object instances, ensuring clear separation and consistent labeling.

Figure 24. ours on SketchyScene inputs (first row), clipasso two styles (second row), sketchagent(last row)

11