

Instance Segmentation of Scene Sketches Using Natural Image Priors

MIA TANG, Stanford University, USA

Yael Vinker, MIT CSAIL, USA

CHUAN YAN, Stanford University, USA

LVMIN ZHANG, Stanford University, USA

MANEESH AGRAWALA, Stanford University, USA



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

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 InkLayer, 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, InkScenes, featuring sketches with diverse brush strokes and varying levels of detail. We use this dataset to demonstrate the robustness

Authors' addresses: Mia Tang, miatang@cs.stanford.edu, Stanford University, Stanford, CA, USA; Yael Vinker, yaelvi116@gmail.com, MIT CSAIL, Boston, USA; Chuan Yan, chuanyan@stanford.edu, Stanford University, Stanford, CA, USA; Lvmmin Zhang, lvmin@stanford.edu, Stanford University, Stanford, CA, USA; Maneesh Agrawala, maneesh@cs.stanford.edu, Stanford University, Stanford, CA, USA.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

SIGGRAPH Conference Papers '25, August 10–14, 2025, Vancouver, BC, Canada

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 979-8-4007-1540-2/2025/08
<https://doi.org/10.1145/3721238.3730606>

of our approach. Code and data for this paper are released at project page: <https://inklayer.github.io>.

CCS Concepts: • Computing methodologies → Image segmentation; Non-photorealistic rendering; Scene understanding.

Additional Key Words and Phrases: Sketch Segmentation, Sketch Understanding, Sketch Editing

ACM Reference Format:

Mia Tang, Yael Vinker, Chuan Yan, Lvmmin Zhang, and Maneesh Agrawala. 2025. Instance Segmentation of Scene Sketches Using Natural Image Priors. In *Special Interest Group on Computer Graphics and Interactive Techniques Conference Conference Papers (SIGGRAPH Conference Papers '25)*, August 10–14, 2025, Vancouver, BC, Canada. ACM, New York, NY, USA, 29 pages. <https://doi.org/10.1145/3721238.3730606>

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 [Chowdhury et al. 2022; Sangkloy et al. 2016; Zou et al. 2018]. 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 InkLayer, 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 segmentation map to divide the sketch into sorted layers to support effective sketch editing.

Our method builds upon Grounded SAM [Ren et al. 2024], 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 [Liu et al. 2023b] for object detection and Segment Anything (SAM) [Kirillov et al. 2023] for mask generation. We analyze the performance of these models on sketch, 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 [2018]. Our fine-tuning technique achieves a substantial improvement in Grounding DINO’s detection performance on sketches, with Average Precision increasing from 26% to 74%.

For object segmentation, we apply SAM [Kirillov et al. 2023] in the sketch domain, using detected object regions from our finetuned 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 [2022] 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 InkScenes, a synthetic annotated dataset of sketched scenes 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 [Zou et al. 2018], 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 [Krishna et al. 2017], which provides annotated scenes with object variety. Using the InstantStyle method [Wang et al. 2024], we generate expressive, natural-looking sketches spanning 72 categories, extending SketchyScene’s original 45 categories by 53 new categories. Our dataset contains 20,542 annotated scene sketches in total, and is highly extensible. Our evaluations demonstrate that InkLayer 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 [Eitz et al. 2012; Ge et al. 2020; Huang et al. 2014; Li et al. 2018; Wu et al. 2018] to train a segmentation model, and use various network architectures, including CNNs [Wang et al. 2020; Zhu et al. 2018], RNNs [Kaiyrbekov and Sezgin 2020; Qi and Tan 2019; Wu et al. 2018], Graph Neural Networks [Yang et al. 2021; Zheng et al. 2022], Transformers [Wang and Li 2024; Zheng et al. 2023], and more specific techniques such as deformation networks [Qi et al. 2021] and CRFs [Schneider and Tuytelaars 2016]. 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 [Li et al. 2022, 2018, 2019] to achieve class-agnostic segmentation. However 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. [2015] 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. [2018] 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 [He et al. 2018]. 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. [2022] introduced SKY-Scene and TUB-Scene by replacing its object components with sketches from the Sketchy [2016] and TU-Berlin [2012] 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 [Zhang et al. 2022a] 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 [Yang et al. 2023] 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. [2024] finetune the CLIP image encoder [Radford et al. 2021] on the FS-COCO [Chowdhury et al. 2022] 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.

2.3 Image Segmentation

The task of image segmentation have been widely explored [Bolya et al. 2019; Cheng et al. 2022; He et al. 2018; Wang et al. 2021]. The advent of vision-language models [Liu et al. 2023a; Radford et al.

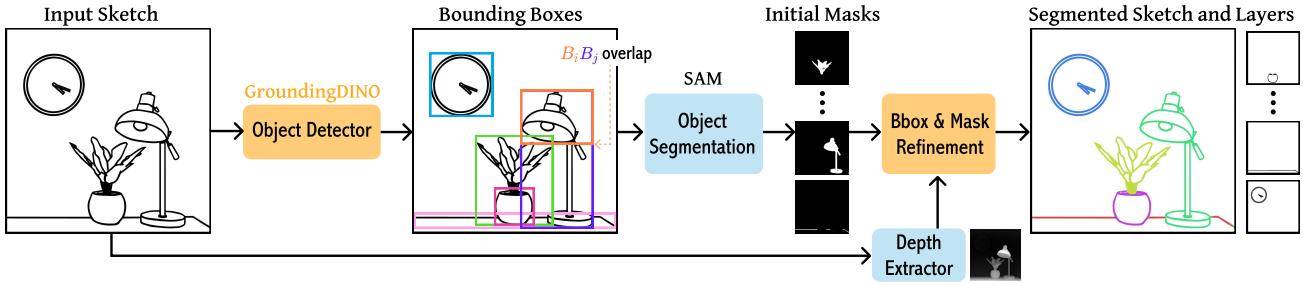


Fig. 2. Overview of the sketch segmentation pipeline. Given an input sketch image, our framework first detects bounding boxes using a customized Grounding DINO to obtain region proposals, and then perform segmentation with SAM models. The localization and segmentation are refined by incorporating the depth features. The result segmentation can be viewed as a layered decomposition of object components in the original sketch.

2021; Xiao et al. 2023] has led to numerous object detection and segmentation methods with impressive generalization capabilities [Kirillov et al. 2023; Minderer et al. 2022; Ren et al. 2024; Zhang et al. 2022b]. Grounding DINO [Liu et al. 2023b] is a state-of-the-art object detection model trained on over 10 million images. It builds on top of DINO [Caron et al. 2021], a strong vision encoder, with effective grounding module that fuses visual and textual information, enabling open-vocabulary detection of unseen objects. Segment Anything (SAM) [Kirillov et al. 2023] 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 [Ren et al. 2024], 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) [Kirillov et al. 2023] model. Next, we perform a refinement stage that leverages scene depth information to assign the final segmentation. This stage also employs a pre-trained inpainting model [von Platen et al. 2022] to produce scene layers.

3.1 Sketch-Aware Object Detection

Grounding DINO [Liu et al. 2023b] 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 [Zou et al. 2018]. 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 [Liu et al. 2021] 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 [Kirillov et al. 2023] to extract masks for 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 have 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

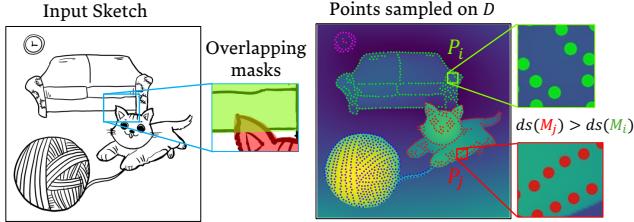


Fig. 3. Resolving ambiguities in overlapping regions. The depth map D is sampled along sketch pixels at evenly spaced points, and the sampled points are grouped by their corresponding object (e.g., P_i corresponds to the i 'th object). Each object is assigned a depth score based on the majority of depth values from the sampled points. Ambiguous pixels are then assigned to the mask with the highest depth score, prioritizing foreground objects.

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 S : $M_i * S$, $M_j * S$, and compute the IoU of these regions to define an “overlapping” score between two objects i, j :

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

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 DepthAnythingV2 [Yang et al. 2024] to extract the depth map D of the input sketch. We then sample D along the sketch pixels at equally spaced points $P = \{p_1, p_2, \dots, p_n\}$, 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:

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

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 [Vincent and Soille 1991] refinement, propagating existing mask labels to 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 pretrained SDXL inpainting model [von Platen et al. 2022]. 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 4). We then identify the group of masks that intersect with M_i : $\mathcal{H}(M_i) = \{M_j | M_j \cap M_i \neq \emptyset\}$. Finally,

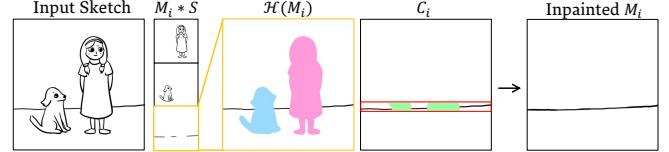


Fig. 4. Layer completion. Object layers are isolated and inpainted using a pretrained SDXL model. The inpainting mask for each object is defined by intersecting overlapping masks with the object’s bounding box.

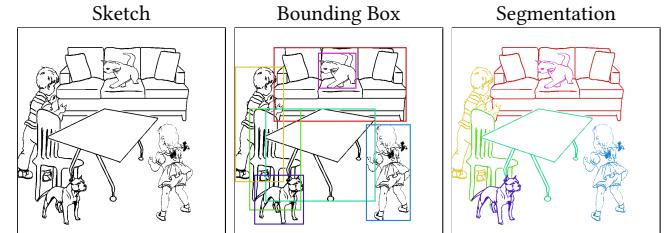


Fig. 5. SketchyScene dataset provides ground truth object bounding boxes and pixel-level instance segmentation masks for scene layouts.

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. 4), 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.

4 SCENE SKETCH SEGMENTATION BENCHMARK

To evaluate our performance across a diverse set of sketches, we construct a synthetic annotated scene-sketch dataset: InkScenes. 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 [Zou et al. 2018] 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

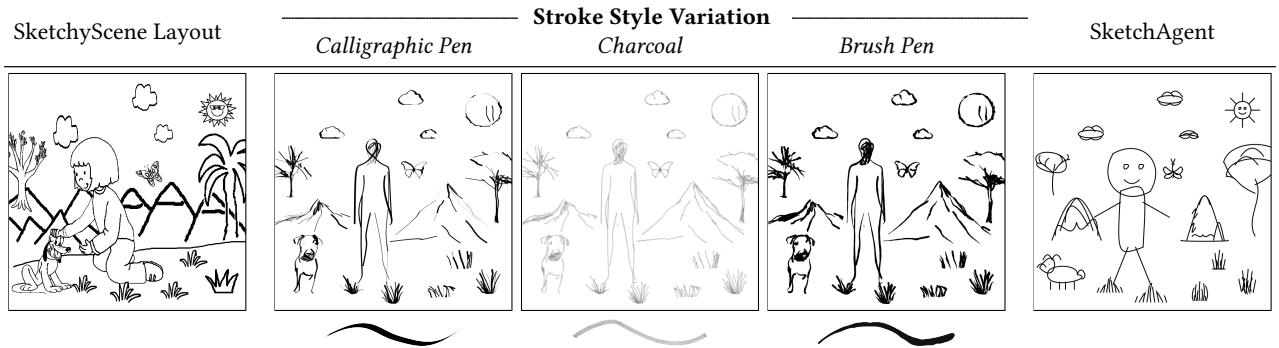


Fig. 6. Samples from our InkScenes dataset. We augment the SketchyScene dataset by generating vector sketches with varied drawing styles based on SketchyScene’s scene layouts. Stroke style variation is introduced by re-rendering the scenes with three different brush styles. Additionally, we create a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



Fig. 7. Illustration of our InkScenes dataset. The input images are sourced from the Visual Genome dataset [Krishna et al. 2017], which we filter to a subset of scenes containing 5 to 10 object instances. We generate the corresponding sketches with InstantStyle [Wang et al. 2024].

styles and stroke variations. Specifically, we incorporate recent object sketching methods that introduce significantly different sketch appearances compared to SketchyScene. These include CLIPasso [Vinker et al. 2022], which transforms images of individual objects into sketches with relatively high image fidelity, and SketchAgent [Vinker et al. 2024], 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 using 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 [Krishna et al. 2017], a large-scale dataset containing diverse and richly

annotated images containing over 33,877 distinct object categories. We use InstantStyle [Wang et al. 2024], 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 InstantStyle subset of the InkScenes dataset includes a total of 72 categories, featuring 53 novel object classes not present in SketchyScene. A few examples of the resulting dataset are shown in Figure 7.

5 RESULTS

Figures 1, 8, 10 and 11 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 [Zou et al. 2018], Sketch-Seger [Yang et al. 2023], and the method proposed by Bourouis *et al.* [2024]. We include additional baselines: Grounding DINO [Liu et al. 2023b] applied directly to sketches, automatic labeling with the Recognize-Anything Model (RAM) [Zhang et al. 2023], and the Automatic SAM pipeline, which samples a 32×32 uniform grid of prompts and selects the top 50 masks to form disjoint regions. Our evaluation dataset consists of 8076 samples: 1,113 test samples from the SketchyScene dataset, 330 samples from the Zhang *et al.* dataset, 1,113 samples from each of the CLIPasso brush styles: Base, Calligraphic Pen, Charcoal, and Brush Pen, 1,113 samples from the SketchAgent, and 1068 samples from the InstantStyle sketch

Table 1. Quantitative comparisons for object detection. We report IoU, AR, and AP metrics across eight 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.

| Metric \ Dataset | 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 |
| Zhang <i>et al.</i> | 0.33 | 0.53 | 0.64 | 0.18 | 0.48 | 0.69 | 0.14 | 0.38 | 0.62 | 0.31 | 0.57 | 0.62 | 0.06 | 0.41 | 0.72 |
| 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 | 0.40 | 0.32 | 0.70 | 0.28 | 0.31 | 0.79 | 0.23 | 0.26 | 0.74 | 0.55 | 0.37 | 0.85 | 0.11 | 0.29 | 0.79 |
| | ± 0.10 | ± 0.12 | ± 0.11 | ± 0.09 | ± 0.10 | ± 0.09 | ± 0.09 | ± 0.08 | ± 0.12 | ± 0.19 | ± 0.13 | ± 0.12 | ± 0.04 | ± 0.09 | ± 0.12 |

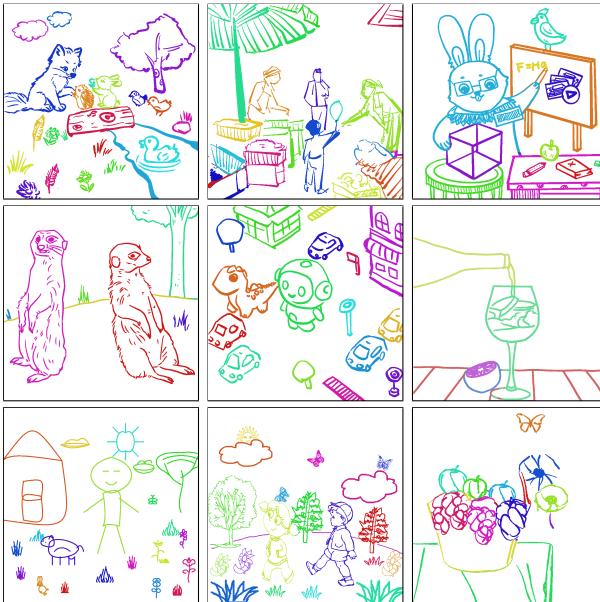


Fig. 8. InkLayer segmentation results. Our method handles sketches of diverse styles and levels of complexity. The first two rows show artist-drawn sketches spanning a range of scenarios, including wildlife scenes, crowded urban markets, and cartoon characters in various settings.

dataset. 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. Therefore, we used SAM for mask generation based on their detected bounding boxes, which yielded better performance 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.* [2024] and SketchSeger [2023] are 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.

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

| Dataset \ Method | Acc ↑ | | | | IoU ↑ | | | |
|---------------------|--------------|------------|------------|-------------|--------------|------------|------------|-------------|
| | SketchyScene | G-SAM | Auto-SAM | Ours | SketchyScene | G-SAM | Auto-SAM | Ours |
| SketchyScene | 0.79 | 0.54 | 0.24 | 0.92 | 0.72 | 0.26 | 0.17 | 0.88 |
| Zhang <i>et al.</i> | 0.42 | 0.72 | 0.09 | 0.86 | 0.33 | 0.51 | 0.04 | 0.74 |
| SketchAgent | 0.39 | 0.35 | 0.12 | 0.88 | 0.38 | 0.16 | 0.06 | 0.84 |
| C-Base | 0.70 | 0.54 | 0.08 | 0.91 | 0.64 | 0.32 | 0.07 | 0.88 |
| C-Calligraphic | 0.66 | 0.50 | 0.06 | 0.87 | 0.63 | 0.30 | 0.05 | 0.86 |
| C-Charcoal | 0.59 | 0.47 | 0.12 | 0.85 | 0.43 | 0.26 | 0.06 | 0.84 |
| C-BrushPen | 0.66 | 0.51 | 0.05 | 0.89 | 0.54 | 0.29 | 0.05 | 0.85 |
| InstantStyle | 0.43 | 0.65 | 0.25 | 0.70 | 0.32 | 0.44 | 0.16 | 0.78 |
| All | 0.58 | 0.53 | 0.13 | 0.87 | 0.50 | 0.32 | 0.08 | 0.82 |
| | ± 0.15 | ± 0.11 | ± 0.08 | ± 0.04 | ± 0.16 | ± 0.11 | ± 0.05 | ± 0.07 |

Table 3. Quantitative comparisons from image segmentation on the filtered datasets. We report Accuracy and IoU metrics across eight datasets, along with the mean and standard deviation for each method. Since Bourouis *et al.* performs semantic segmentation, and requires an input text prompt, we provide ground truth class labels as prompts to generate segmentations, and use 0.01 for confidence threshold to ensure all sketch pixels are segmented.

| Dataset \ Method | Acc ↑ | | | IoU ↑ | | |
|---------------------|------------|-------------|-------------|------------|-------------|-------------|
| | Bourouis | SketchSeger | Ours | Bourouis | SketchSeger | Ours |
| SketchyScene | 0.79 | 0.96 | 0.93 | 0.58 | 0.90 | 0.90 |
| Zhang <i>et al.</i> | 0.69 | 0.65 | 0.86 | 0.50 | 0.22 | 0.75 |
| SketchAgent | 0.71 | 0.69 | 0.93 | 0.52 | 0.48 | 0.88 |
| C-Base | 0.83 | 0.86 | 0.94 | 0.66 | 0.73 | 0.92 |
| C-Calligraphic | 0.79 | 0.76 | 0.88 | 0.61 | 0.61 | 0.88 |
| C-Charcoal | 0.80 | 0.80 | 0.85 | 0.61 | 0.57 | 0.85 |
| C-BrushPen | 0.81 | 0.79 | 0.91 | 0.63 | 0.59 | 0.88 |
| InstantStyle | 0.76 | 0.70 | 0.80 | 0.51 | 0.45 | 0.75 |
| All | 0.77 | 0.78 | 0.89 | 0.58 | 0.57 | 0.85 |
| | ± 0.05 | ± 0.10 | ± 0.05 | ± 0.06 | ± 0.20 | ± 0.07 |

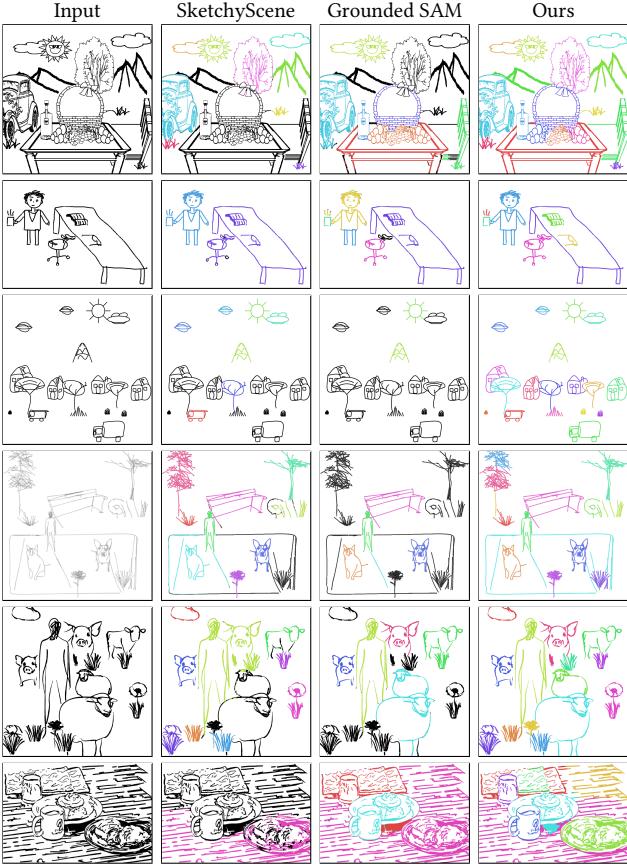


Fig. 9. Qualitative comparison of instance segmentation methods. Each row corresponds to a different sketch dataset. Black pixels indicate regions where segmentation was not applied. Our method effectively segments sketch pixels into distinct instances, outperforming alternative approaches.

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 Table 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 0.79, 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 without adaptation, despite its strong performance on natural images. This is evident from the large



Fig. 10. InkLayer results on artist-drawn and freehand sketches. Our method accurately segments artist-drawn sketches in the first three rows and effectively handles quick, novice freehand sketches from the Zhang *et al.* [2018] and FSCOCO-Seg [2024] datasets in the bottom two rows.

margin in scores between our method and Grounding DINO, seeing an increase of 38% in IoU, 48% increase in AR, and 48% 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 Table 2, while the results for Bourouis *et al.* [2024] and SketchSeger [2023] are shown separately in Table 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. While SketchSeger performs well on SketchyScene dataset, its performance degrades significantly on the symbolic style SketchAgent and struggles to generalize to novel categories included in Zhang *et al.* and InstantStyle styles. Meanwhile, our method demonstrates robustness to various styles, especially excelling on the challenging Zhang *et al.* and SketchAgent style compared to other methods.



Fig. 11. **InkLayer segmentation results on synthetic sketches.** Our method is able to precisely segment SketchyScene dataset (first row) and our InkScenes dataset: CLIPasso variants (second and third rows), abstract SketchAgent sketches (fourth row), and detailed InstantStyle sketches (last row).

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 12a). 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 12b). In particular, it may miss sketch object boundaries or produce grid-like masks. 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 InkLayer, 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 scene-level annotated sketch dataset, InkScenes, encompassing a wide range of

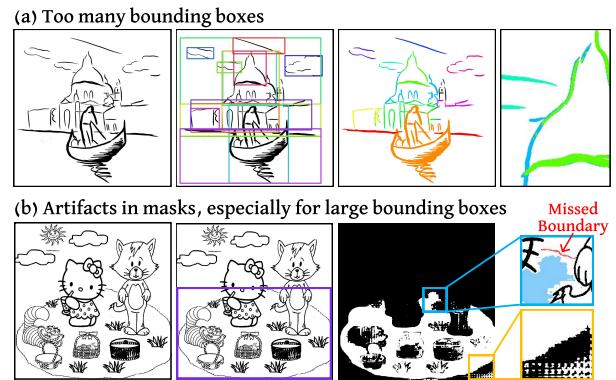


Fig. 12. **Examples showcasing InkLayer’s limitations.** (a) The bounding box filtering process can still retain undesired boxes, leading to artifacts in final segmentation. (b) SAM masks are generally reasonable but can produce artifacts such as missing object boundaries or producing grid-like regions.

object categories and significant variations in drawing styles. Our experiments demonstrate that InkLayer is robust to these variations, achieving consistent performance across diverse datasets.

ACKNOWLEDGMENTS

This work was partially supported by the Brown Institute for Media Innovation at Stanford. Yael Vinker was supported in part by IBM Agreement No. W1771646 and Hyundai Motor Company R&D Center Agmt Dtd 2/22/23.

- Daniel Bolya, Chong Zhou, Fanyi Xiao, and Yong Jae Lee. 2019. YOLACT: Real-Time Instance Segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.
- Ahmed Bourouis, Judith E Fan, and Yulia Gryaditskaya. 2024. Open Vocabulary Semantic Scene Sketch Understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4176–4186.
- Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. 2021. Emerging Properties in Self-Supervised Vision Transformers. In *ICCV*.
- Bowen Cheng, Ishan Misra, Alexander G. Schwing, Alexander Kirillov, and Rohit Girdhar. 2022. Masked-attention Mask Transformer for Universal Image Segmentation. *CVPR*.
- Pinaki Nath Chowdhury, Aneeshan Sain, Ayan Kumar Bhunia, Tao Xiang, Yulia Gryaditskaya, and Yi-Zhe Song. 2022. FS-COCO: Towards Understanding of Freehand Sketches of Common Objects in Context. In *ECCV*.
- Mathias Eitz, James Hays, and Marc Alexa. 2012. How Do Humans Sketch Objects? *ACM Trans. Graph. (Proc. SIGGRAPH)* 31, 4 (2012), 44:1–44:10.
- Ce Ge, Haifeng Sun, Yi-Zhe Song, Zhanyu Ma, and Jianxin Liao. 2022. Exploring Local Detail Perception for Scene Sketch Semantic Segmentation. *IEEE Transactions on Image Processing* 31 (2022), 1447–1461. <https://doi.org/10.1109/TIP.2022.3142511>
- Songwei Ge, Vedanuj Goswami, C. Lawrence Zitnick, and Devi Parikh. 2020. Creative Sketch Generation. arXiv:2011.10039 [cs.CV]
- Yulia Gryaditskaya, Mark Pesteyn, Jan Willem Hoftijzer, Sylvia Pont, Frédéric Durand, and Adrien Bousseau. 2019. OpenSketch: A Richly-Annotated Dataset of Product Design Sketches. *ACM Transactions on Graphics (Proc. SIGGRAPH Asia)* 38 (11 2019).
- Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick. 2018. Mask R-CNN. arXiv:1703.06870 [cs.CV] <https://arxiv.org/abs/1703.06870>
- Zhe Huang, Hongbo Fu, and Rynson W. H. Lau. 2014. Data-driven segmentation and labeling of freehand sketches. *ACM Trans. Graph.* 33, 6, Article 175 (Nov. 2014), 10 pages. <https://doi.org/10.1145/2661229.2661280>
- Kurmanbek Kaiyrbekov and Metin Sezgin. 2020. Deep Stroke-Based Sketched Symbol Reconstruction and Segmentation. *IEEE Computer Graphics and Applications* 40, 1 (2020), 112–126. <https://doi.org/10.1109/MCG.2019.2943333>
- 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. 2023. Segment Anything. arXiv:2304.02643 (2023).
- 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. 2017. Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations. *Int. J. Comput. Vision* 123, 1 (May 2017), 32–73. <https://doi.org/10.1007/s11263-016-0981-7>
- Changjian Li, Hao Pan, Adrien Bousseau, and Niloy J. Mitra. 2022. Free2CAD: Parsing Freehand Drawings into CAD Commands. *ACM Trans. Graph. (Proceedings of SIGGRAPH 2022)* 41, 4 (2022), 93:1–93:16. <https://doi.org/10.1145/3528223.3530133>
- Ke Li, Kaiyue Pang, Jifei Song, Yi-Zhe Song, Tao Xiang, Timothy M. Hospedales, and Honggang Zhang. 2018. Universal Perceptual Grouping. arXiv:1808.02312 [cs.CV] <https://arxiv.org/abs/1808.02312>
- Ke Li, Kaiyue Pang, Yi-Zhe Song, Tao Xiang, Timothy M. Hospedales, and Honggang Zhang. 2019. Toward Deep Universal Sketch Perceptual Grouper. *IEEE Transactions on Image Processing* 28 (2019), 3219–3231. <https://api.semanticscholar.org/CorpusID:73430225>
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023a. Visual Instruction Tuning. In *NeurIPS*.
- Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. 2023b. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499* (2023).
- Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. 2021. Swin Transformer: Hierarchical Vision Transformer using Shifted Windows. In *ICCV*.
- 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. 2022. Simple Open-Vocabulary Object Detection with Vision Transformers. arXiv:2205.06230 [cs.CV] <https://arxiv.org/abs/2205.06230>
- Anran Qi, Yulia Gryaditskaya, Tao Xiang, and Yi-Zhe Song. 2021. One Sketch for All: One-Shot Personalized Sketch Segmentation. *CoRR* abs/2112.10838 (2021). arXiv:2112.10838 <https://arxiv.org/abs/2112.10838>
- Yonggang Qi, Yi-Zhe Song, Tao Xiang, Honggang Zhang, Timothy Hospedales, Yi Li, and Jun Guo. 2015. Making better use of edges via perceptual grouping. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 1856–1865. <https://doi.org/10.1109/CVPR.2015.7298795>
- Yonggang Qi and Zheng-Hua Tan. 2019. SketchSegNet+: An End-to-End Learning of RNN for Multi-Class Sketch Semantic Segmentation. *IEEE Access* 7 (2019), 102717–102726. <https://doi.org/10.1109/ACCESS.2019.2929804>
- 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. 2021. Learning Transferable Visual Models From Natural Language Supervision. arXiv:2103.00020 [cs.CV] <https://arxiv.org/abs/2103.00020>
- 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. 2024. Grounded SAM: Assembling Open-World Models for Diverse Visual Tasks. arXiv:2401.14159 [cs.CV]
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2021. High-Resolution Image Synthesis with Latent Diffusion Models. arXiv:2112.10752 [cs.CV]
- Patsorn Sangkloy, Nathan Burnell, Cusuh Ham, and James Hays. 2016. The sketchy database: learning to retrieve badly drawn bunnies. *ACM Trans. Graph.* 35, 4, Article 119 (July 2016), 12 pages. <https://doi.org/10.1145/2897824.2925954>
- Rosália G. Schneider and Tinne Tuytelaars. 2016. Example-Based Sketch Segmentation and Labeling Using CRFs. *ACM Trans. Graph.* 35, 5, Article 151 (July 2016), 9 pages. <https://doi.org/10.1145/2898351>
- L. Vincent and P. Soille. 1991. Watersheds in digital spaces: an efficient algorithm based on immersion simulations. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 13, 6 (1991), 583–598. <https://doi.org/10.1109/34.87344>
- Yael Vinker, Ehsan Pajouheshgar, Jessica Y. Bo, Roman Christian Bachmann, Amit Haim Bernano, Daniel Cohen-Or, Amir Zamir, and Ariel Shamir. 2022. CLIPasso: Semantically-Aware Object Sketching. *ACM Trans. Graph.* 41, 4, Article 86 (Jul 2022), 11 pages. <https://doi.org/10.1145/3528223.3530068>
- Yael Vinker, Tamar Rott Shaham, Kristine Zheng, Alex Zhao, Judith E Fan, and Antonio Torralba. 2024. SketchAgent: Language-Driven Sequential Sketch Generation. arXiv:2411.17673 [cs.CV] <https://arxiv.org/abs/2411.17673>
- 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. 2022. Diffusers: State-of-the-art diffusion models. <https://github.com/huggingface/diffusers>.
- Fei Wang, Shujin Lin, Hanhui Li, Hefeng Wu, Tie Cai, Xiaonan Luo, and Ruomei Wang. 2020. Multi-column point-CNN for sketch segmentation. *Neurocomputing* 392 (2020), 50–59. <https://doi.org/10.1016/j.neucom.2019.12.117>
- Haofan Wang, Matteo Spinelli, Qixun Wang, Xu Bai, Zekui Qin, and Anthony Chen. 2024. InstantStyle: Free Lunch towards Style-Preserving in Text-to-Image Generation. *ArXiv* abs/2404.02733 (2024). <https://api.semanticscholar.org/CorpusID:268876474>
- Jiawei Wang and Changjian Li. 2024. ContextSeg: Sketch Semantic Segmentation by Querying the Context with Attention. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3679–3688.
- Xinlong Wang, Rufeng Zhang, Chunhua Shen, Tao Kong, and Lei Li. 2021. SOLO: A Simple Framework for Instance Segmentation. *IEEE T. Pattern Analysis and Machine Intelligence (TPAMI)* (2021).
- Xingyu Wu, Yonggang Qi, Jun Liu, and Jie Yang. 2018. Sketchsegnet: A Rnn Model for Labeling Sketch Strokes. *2018 IEEE 28th International Workshop on Machine Learning for Signal Processing (MLSP)* (2018), 1–6. <https://api.semanticscholar.org/CorpusID:53234925>
- Bin Xiao, Haiping Wu, Weijian Xu, Xiyang Dai, Houdong Hu, Yumao Lu, Michael Zeng, Ce Liu, and Lu Yuan. 2023. Florence-2: Advancing a unified representation for a variety of vision tasks. *arXiv preprint arXiv:2311.06242* (2023).
- Jie Yang, Aihua Ke, Yaoxiang Yu, and Bo Cai. 2023. Scene sketch semantic segmentation with hierarchical Transformer. *Knowledge-Based Systems* 280 (2023), 110962.
- Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. 2024. Depth Anything V2. arXiv:2406.09414 (2024).
- Lumin Yang, Jiajie Zhuang, Hongbo Fu, Xiangzhi Wei, Kun Zhou, and Youyi Zheng. 2021. SketchGNN: Semantic Sketch Segmentation with Graph Neural Networks. *ACM Trans. Graph.* 40, 3, Article 28 (Aug. 2021), 13 pages. <https://doi.org/10.1145/3450284>
- Haotian* Zhang, Pengchuan* Zhang, Xiaowei Hu, Yen-Chun Chen, Liunian Harold Li, Xiyang Dai, Lijuan Wang, Lu Yuan, Jenq-Neng Hwang, and Jianfeng Gao. 2022b. GLIPv2: Unifying Localization and Vision-Language Understanding. *arXiv preprint arXiv:2206.05836* (2022).
- Jianhui Zhang, Yilai Chen, Lei Li, Hongbo Fu, and Chiew-Lan Tai. 2018. Context-based sketch classification. In *Proceedings of the Joint Symposium on Computational Aesthetics and Sketch-Based Interfaces and Modeling and Non-Photorealistic Animation and Rendering* (Victoria, British Columbia, Canada) (*Expressive ’18*). Association for Computing Machinery, New York, NY, USA, Article 3, 10 pages. <https://doi.org/10.1145/3229147.3229154>
- Youcai Zhang, Xinyu Huang, Jinyu Ma, Zhaoyang Li, Zhaochuan Luo, Yanchun Xie, Yuzhuo Qin, Tong Luo, Yaqian Li, Shilong Liu, et al. 2023. Recognize Anything: A Strong Image Tagging Model. *arXiv preprint arXiv:2306.03514* (2023).
- Zhengming Zhang, Xiaoming Deng, Jinyao Li, Yukun Lai, Cuixia Ma, Yongjin Liu, and Hongan Wang. 2022a. Stroke-based semantic segmentation for scene-level free-hand sketches. *Vis. Comput.* 39, 12 (Dec. 2022), 6309–6321. <https://doi.org/10.1007/s00371-022-02731-8>
- Xiaoxiao Zheng, Jiyang Xie, Aneeshan Sain, Zhanyu Ma, Yi-Zhe Song, and Jun Guo. 2022. ENDE-GNN: An Encoder-decoder GNN Framework for Sketch Semantic

- Segmentation. In *2022 IEEE International Conference on Visual Communications and Image Processing (VCIP)*. 1–5. <https://doi.org/10.1109/VCIP56404.2022.10008880>
- Xixiao Zheng, Jiyang Xie, Aneeshan Sain, Yi-Zhe Song, and Zhanyu Ma. 2023. Sketch-Segformer: Transformer-Based Segmentation for Figurative and Creative Sketches. *IEEE Transactions on Image Processing* 32 (2023), 4595–4609. <https://doi.org/10.1109/TIP.2023.3302521>
- Xianyi Zhu, Yi Xiao, and Yan Zheng. 2018. Part-Level Sketch Segmentation and Labeling Using Dual-CNN. In *International Conference on Neural Information Processing*. <https://api.semanticscholar.org/CorpusID:54435773>
- Changqing Zou, Qian Yu, Ruofei Du, Haoran Mo, Yi-Zhe Song, Tao Xiang, Chengying Gao, Baoguan Chen, and Hao Zhang. 2018. SketchyScene: Richly-Annotated Scene Sketches. In *ECCV*. Springer International Publishing, 438–454. https://doi.org/10.1007/978-3-030-01267-0_26

Instance Segmentation of Scene Sketches Using Natural Image Priors

Supplementary Material

Mia Tang¹ Yael Vinker² Chuan Yan¹ Lvmmin Zhang¹ Maneesh Agrawala¹

¹Stanford University ²MIT CSAIL

Project page: <https://inklayer.github.io>

CONTENTS

| | |
|---------------------------------------------|----|
| Contents | 11 |
| A Sketch Editing Interface | 11 |
| B Refinement Module Ablation | 11 |
| C Depth Maps Accuracy | 11 |
| C.1 Qualitative Evaluation of Depth Maps | 13 |
| C.2 Quantitative Evaluation of Depth Maps | 13 |
| D InkScenes Dataset Generation Details | 13 |
| D.1 Generating Vector Scene Sketches | 13 |
| D.2 Generating Sketches from Natural Images | 13 |
| E InkScenes Dataset Statistics | 13 |
| F Performance on Human-Drawn Sketches | 15 |
| F.1 Performance on FSCOCO-Seg | 15 |
| F.2 Performance on Product Design Sketches | 15 |
| F.3 Performance on Zhang <i>et al.</i> | 16 |
| G Additional Qualitative Comparisons | 16 |

A SKETCH EDITING INTERFACE



Fig. 13. **InkLayer application.** Interactive interface for sketch editing, powered by our instance segmentation and layer completion algorithm.

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)

| Dataset | Metric | | Acc ↑ | | IoU ↑ | |
|---------------------|----------|-------------|----------|-------------|----------|---------|
| | w/o Ref. | w/ Ref. | w/o Ref. | w/ Ref. | w/o Ref. | w/ Ref. |
| SketchyScene | 0.85 | 0.92 | 0.79 | 0.88 | | |
| Zhang <i>et al.</i> | 0.79 | 0.86 | 0.67 | 0.74 | | |
| SketchAgent | 0.82 | 0.88 | 0.76 | 0.84 | | |
| CLIPasso base | 0.86 | 0.91 | 0.81 | 0.88 | | |
| CLIPasso 01 | 0.82 | 0.87 | 0.79 | 0.86 | | |
| CLIPasso 04 | 0.85 | 0.85 | 0.82 | 0.84 | | |
| CLIPasso 11 | 0.82 | 0.89 | 0.77 | 0.85 | | |
| InstantStyle | 0.73 | 0.78 | 0.60 | 0.70 | | |
| All | 0.82 | 0.87 | 0.75 | 0.82 | | |
| | ± 0.04 | ± 0.04 | ± 0.08 | ± 0.07 | | |

Table 4. Comparison of segmentation performance with and without refinement module. We see consistent improvement in performance across all datasets.

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.

B REFINEMENT MODULE ABLATION

We conduct an ablation study on our depth-guided refinement module across the benchmark datasets, with quantitative results summarized in Table 4. On average, removing the refinement module results in segmentation performance degradation of a 0.05 decrease in accuracy and a 0.07 decrease in IoU, underscoring its importance in achieving precise instance segmentation. Figure 14 presents qualitative comparisons of segmentation results for the same sketch, with and without the refinement module. These examples further illustrate the module’s effectiveness in resolving ambiguities and enhancing segmentation quality.

C DEPTH MAPS ACCURACY

The effectiveness of our refinement module depends on accurately sorting object instances, which in turn relies on the quality of the estimated sketch depth maps. To generate these depth maps, we use DepthAnythingV2 [Yang et al. 2024], a model trained on a diverse set of visual inputs including artworks and sketches. However, evaluating depth map accuracy remains challenging due to the lack of ground-truth depth annotations for our synthetic sketches. To address this, we present both qualitative visualizations and a proxy quantitative evaluation. Specifically, we assess depth map quality on



Fig. 14. Refinement module ablation. We compare segmentation outputs with and without our depth-guided refinement module. The refinement significantly improves mask coherence and reduces visual artifacts, such as duplicated or fragmented object regions. As shown across diverse scene sketches, the refined results exhibit smoother boundaries, better alignment with object structure, and fewer inconsistencies—particularly in complex regions with overlapping or adjacent objects.

InstantStyle sketches, for which we have access to the corresponding natural images. These natural images are assumed to yield more reliable depth estimations, enabling a comparative evaluation of the sketch-derived depth maps.

C.1 Qualitative Evaluation of Depth Maps

We present qualitative results of the estimated depth maps in Figure 15, showcasing a variety of sketches across different datasets and scene types. These examples highlight the ability of DepthAnythingV2 to extract meaningful depth cues from sparse and stylized sketch inputs. Despite the abstract nature of the input sketches—ranging from realistic human figures to symbolic line drawings of animals and buildings—the model produces coherent depth estimates that reflect the relative spatial layout of the scene. These depth cues play a crucial role in our refinement module, helping to disambiguate overlapping regions and enforce depth-aware instance separation.

C.2 Quantitative Evaluation of Depth Maps

For the quantitative evaluation of depth quality, we assess the consistency of object ordering between the predicted depth maps (from sketches) and the ground truth depth maps (from corresponding natural images) using Kendall’s Tau coefficient. Rather than relying on absolute depth metrics such as average relative error, we adopt this rank-based approach to better reflect our method’s reliance on relative depth for instance layering. Since our refinement module operates on object ordering rather than precise depth values, this metric offers a more meaningful evaluation. Our results show an average of 80% agreement in object ordering, demonstrating the effectiveness of sketch-based depth estimation for supporting depth-aware segmentation.

D INKSCENES DATASET GENERATION DETAILS

In this section, we provide additional details on our synthetic dataset InkScenes’s creation process.

D.1 Generating Vector Scene Sketches

We employ CLIPasso [2022] and SketchAgent [2024] 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 [2018], ensuring sufficient variability in our synthetic scenes. For CLIPasso’s image-to-vector conversion, we first generate photorealistic synthetic images using SDXL [2021]. The generation process uses a consistent prompt template: “A *realistic image of a {class_name} with a blank background*”. Figure 16 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 17 illustrates representative examples of the generated object sketches.

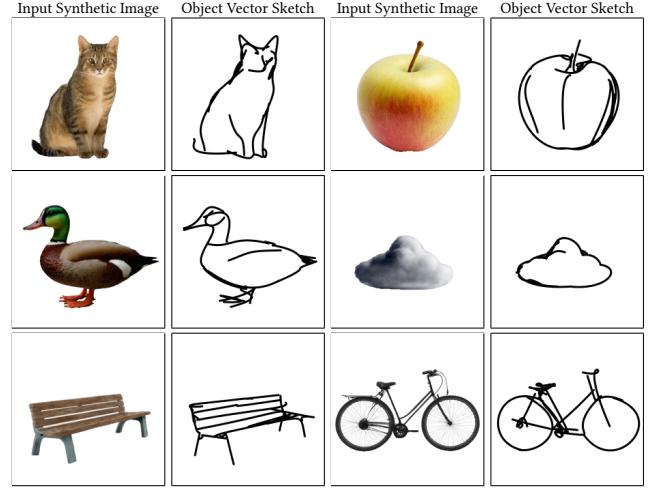


Fig. 16. Examples pairs of input synthetic image and output generated object vector sketch.

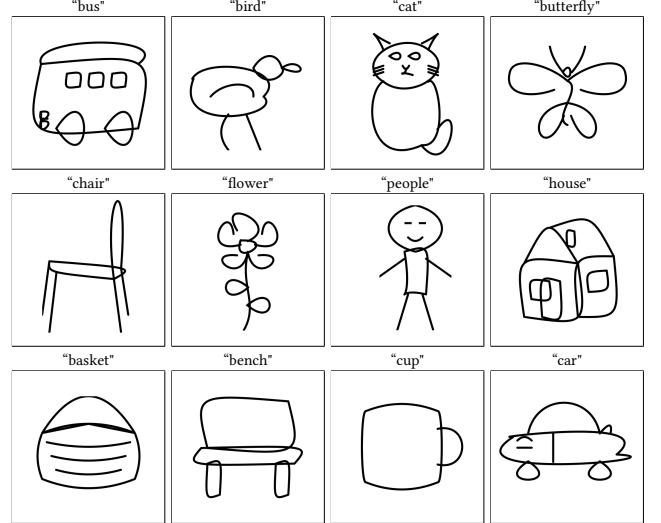


Fig. 17. Examples of object vector sketches generated by SketchAgent.

D.2 Generating Sketches from Natural Images

To expand beyond SketchyScene’s object categories, we employ InstantStyle [2024] to transform a subset of Visual Genome [2017] images into sketches, using a single CLIPasso object sketch as the style reference. Figure 24 showcases a gallery of examples from our InstantStyle-generated scene sketch dataset.

E INKSCENES DATASET STATISTICS

In this section, we present detailed information about our benchmark dataset, *InkScenes*. While the CLIPasso variants and SketchAgent styles of our InkScenes dataset contain the same set of classes and scene complexity as original SketchyScene, our InstantStyle dataset contributes new categories and scene layouts. Figure 18 shows the distribution of categories across the entire dataset, grouped



Fig. 15. Sketch and corresponding DepthAnythingV2 [2024] depth map pairs. We show qualitative results on sketch depth maps of scenes across datasets, demonstrating the reasonable performance of this depth estimator directly on sketches. The model generalizes surprisingly well across a wide spectrum of sketch styles, from highly abstract and symbolic representations to more detailed and realistic ones. Despite the lack of color or shading cues, the predicted depth maps often capture overall scene layout and object ordering effectively.

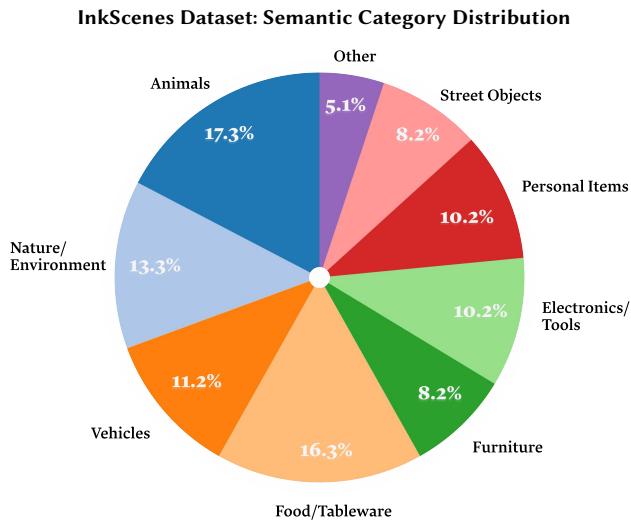


Fig. 18. **Distribution of object categories in the InkScenes dataset, grouped into high-level semantic categories.** We group all 98 unique classes into high-level categories such as animals, nature, vehicles, and household objects.

in high-level semantic categories. Additionally, we show the InstantStyle dataset’s scene complexity in Figure 20, and the class distribution in Figure 19. The full list of categories in our InkScenes dataset is provided in Table 5, grouped into base categories (those included in SketchyScene) and novel categories (unseen and not present in SketchyScene). We omit the “others” category from SketchyScene due to its ambiguity, which makes it unsuitable for generating our new synthetic components.

F PERFORMANCE ON HUMAN-DRAWN SKETCHES

To ensure the practical utility of our method, it is important that it generalizes beyond synthetic datasets to real-world, human-drawn sketches. In the main paper, we evaluate performance on one such dataset introduced by Zhang et al. [2018]. In this section, we extend our evaluation to two additional datasets containing human-drawn sketches, which required more careful consideration in both setup and analysis.

F.1 Performance on FSCOCO-Seg

FSCOCO-Seg [Bourouis et al. 2024] is a freehand scene sketch dataset comprising 500 sketches (originally from FS-COCO [Chowdhury et al. 2022] dataset) with ground-truth semantic segmentation labels, but no instance-level annotations. To evaluate our method on this dataset, we convert our predicted instance segmentation outputs into semantic segmentation maps using the provided ground-truth labels. Our method achieves a pixel accuracy of 0.85 and an IoU of 0.75, demonstrating strong generalization to freehand sketches. Figure 21 shows qualitative instance segmentation results on FSCOCO-Seg, highlighting the method’s ability to segment complete object instances even in abstract and loosely drawn scenes.

| Base Classes (45) | Novel Classes (53) |
|-------------------|--------------------|
| airplane | sandwich |
| apple | refrigerator |
| balloon | couch |
| banana | pizza |
| basket | laptop |
| bee | bed |
| bench | mouse |
| bicycle | toilet |
| bird | orange |
| bottle | toaster |
| bucket | kite |
| bus | cell phone |
| butterfly | cake |
| car | carrot |
| cat | parking meter |
| chair | tv |
| chicken | knife |
| cloud | remote |
| cow | train |
| cup | tie |
| dinnerware | vase |
| dog | potted plant |
| duck | clock |
| fence | sports ball |
| flower | handbag |
| grape | fire hydrant |
| grass | wine glass |
| horse | elephant |
| house | skateboard |
| moon | keyboard |
| mountain | teddy bear |
| people/person | skis |
| picnic rug | backpack |
| pig | spoon |
| rabbit | book |
| road | stop sign |
| sheep | broccoli |
| sofa | zebra |
| star | donut |
| street lamp | sink |
| sun | surfboard |
| table | snowboard |
| tree | motorcycle |
| truck | suitcase |
| umbrella | dining table |
| | boat |
| | fork |
| | microwave |
| | oven |
| | bowl |
| | tennis racket |
| | toothbrush |
| | hot dog |

Table 5. **Full List of Classes in InkScenes Dataset.** The Base Classes are the same as SketchyScene[2018] dataset. The Novel Classes are introduced by our InstantStyle dataset.

F.2 Performance on Product Design Sketches

We do additional qualitative evaluation on product design sketches from OpenSketch dataset [Gryaditskaya et al. 2019] in Fig. 22. Since OpenSketch contains product design sketch of single, individual objects, we manually created scenes by combining multiple object

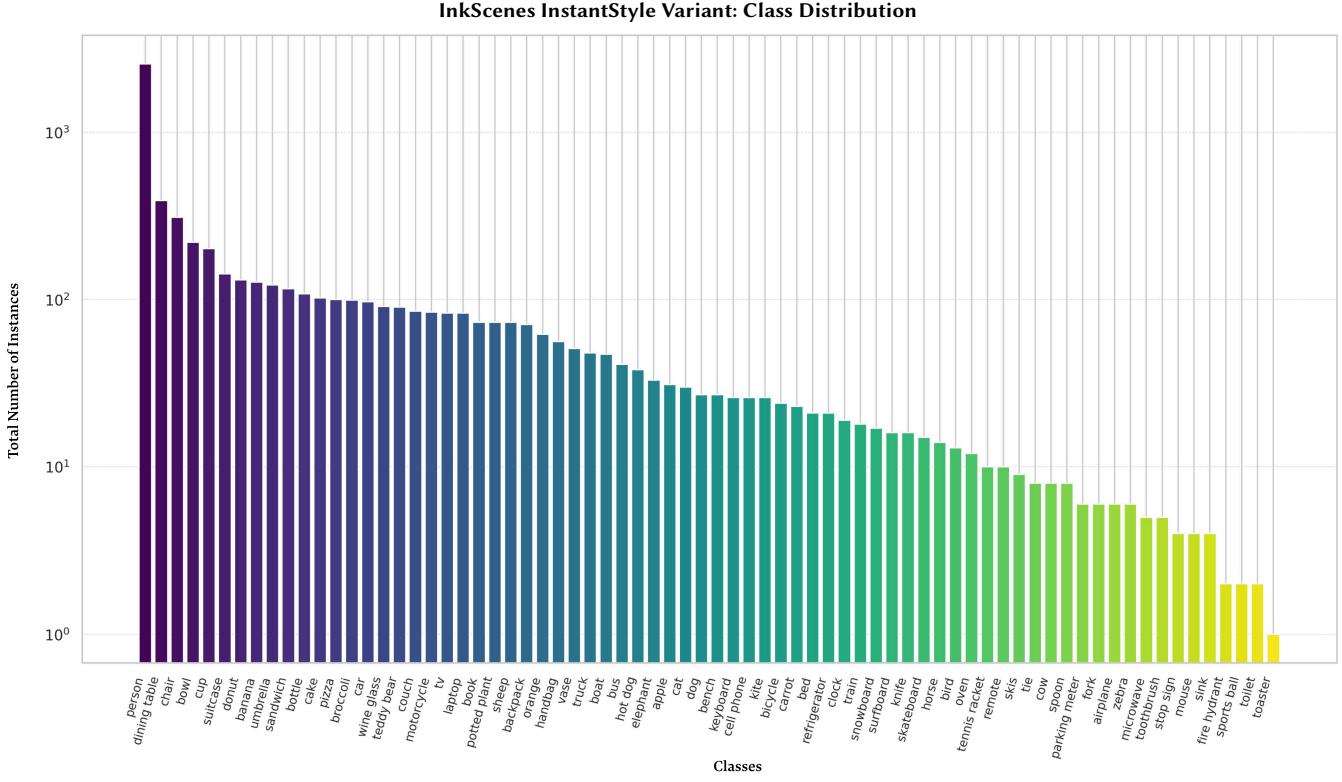


Fig. 19. **Class distribution of the InstantStyle portion of the InkScenes dataset.** By incorporating 53 novel object categories with varied frequencies, our dataset expands beyond prior scene sketch datasets and introduces realistic long-tailed challenges for instance segmentation. This diverse distribution is essential for evaluating model generalization across both common and rare object classes.

InkScenes InstantStyle Variant: Scene Composition Complexity

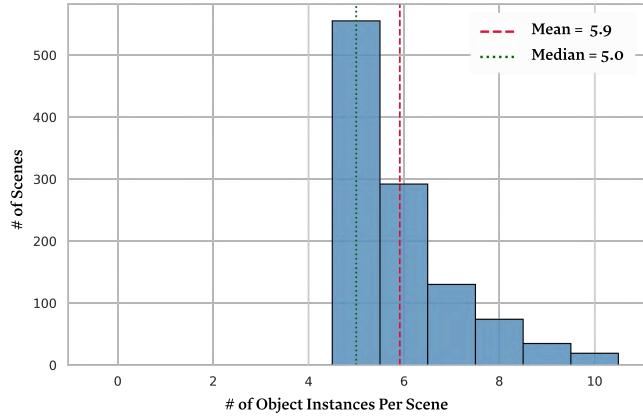


Fig. 20. **Scene-level complexity of the InstantStyle portion of InkScenes.** Each bar counts how many sketches contain a given number of object instances (background excluded). Most scenes cluster around 5–6 objects (median≈5, mean≈6), with a modest long-tail of denser compositions up to ten objects.

sketches. While our method generally succeeds in segmenting individual objects across scenes, we observe limitations such as artifacts

on objects with long, thick construction lines, and segmenting objects apart from their construction lines.

F.3 Performance on Zhang *et al.*

In addition to the numerical evaluations we include in the main paper, we present qualitative results on Zhang *et al.* [2018]’s dataset in Fig 23. This dataset consists of 330 scene sketches drawn by humans, spanning 74 categories with 24 novel categories not included in our InkScenes. Our performance on this dataset further shows our generalization capabilities for unseen categories and scenes.

G ADDITIONAL QUALITATIVE COMPARISONS

We present additional qualitative comparisons between our approach and baseline methods across benchmark scene sketch datasets in Fig. 25 for SketchyScene, Fig. 26 for SketchAgent, Fig. 27 for CLIPasso, Fig. 28 for InstantStyle, to accompany our numerical evaluations included in the paper. To compare with semantic segmentation method, Bourouis *et al.* [2024] and SketchSeger [Yang *et al.* 2023], we created a filtered version of all eight 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. 29 for filtered SketchyScene dataset, Fig. 30 for filtered SketchAgent dataset, and Fig. 31 for filtered CLIPasso dataset, Fig. 32 for filtered Zhang *et al.* dataset, to accompany our qualitative results shown in the paper.



Fig. 21. FSCOCO-Seg Performance. We qualitatively evaluate our method on scenes from the FSCOCO-Seg dataset. Overall, the segmentations are visually accurate and align well with object boundaries in diverse, hand-drawn scenes. While our method performs well across a variety of object categories and layouts, we observe some artifacts in regions with heavy occlusion, where foreground and background objects may blend or be incompletely separated.

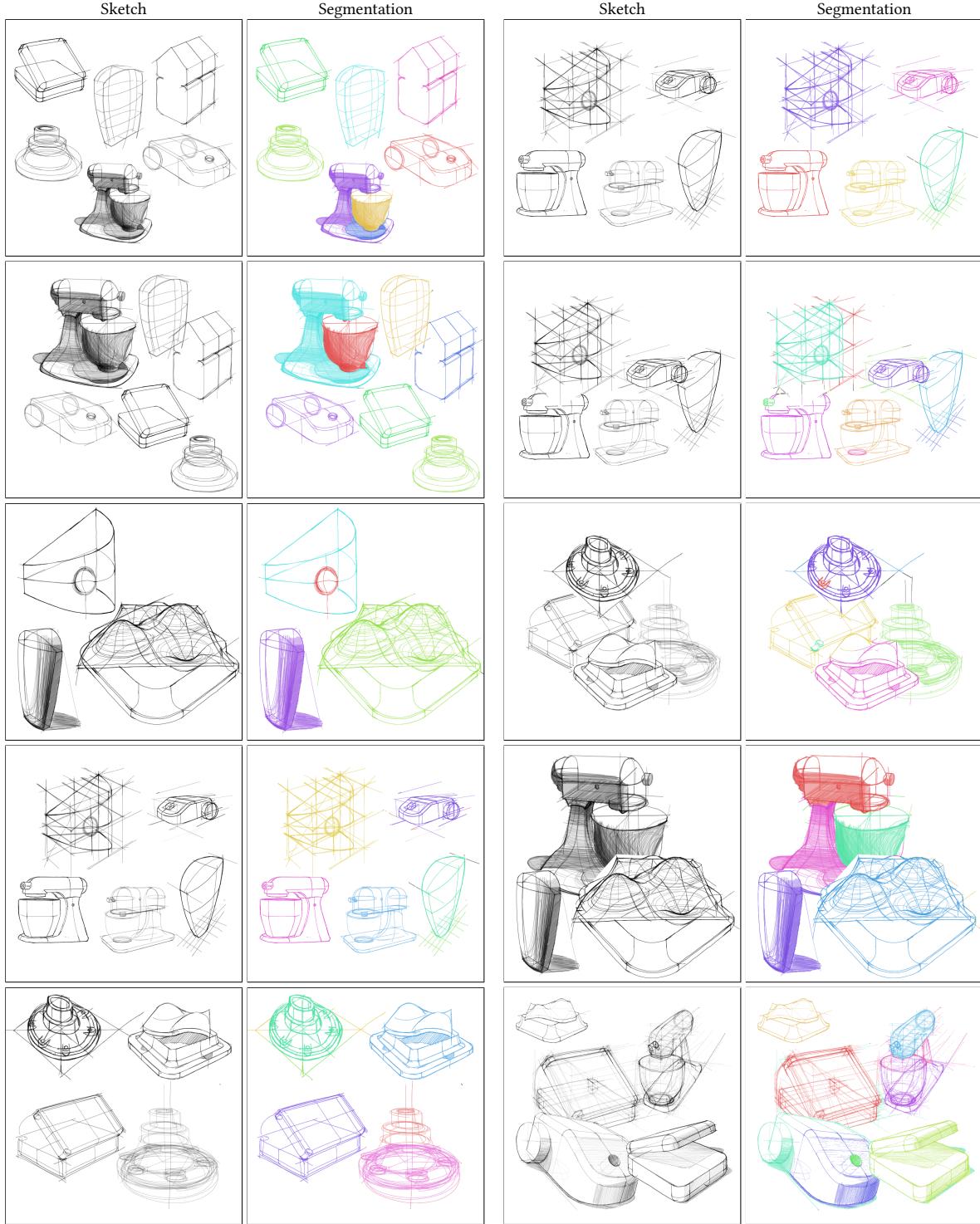


Fig. 22. OpenSketch Performance. We qualitatively evaluate our method on manually composed scenes from the OpenSketch dataset. For clearer occlusion handling, we manually added coarse white backgrounds behind foreground objects. Our method segments individual objects reasonably well, especially when sketches are not heavily cluttered with construction lines. For example, segmentation struggles with the background object in the second example of row 2, which contains dense construction lines.

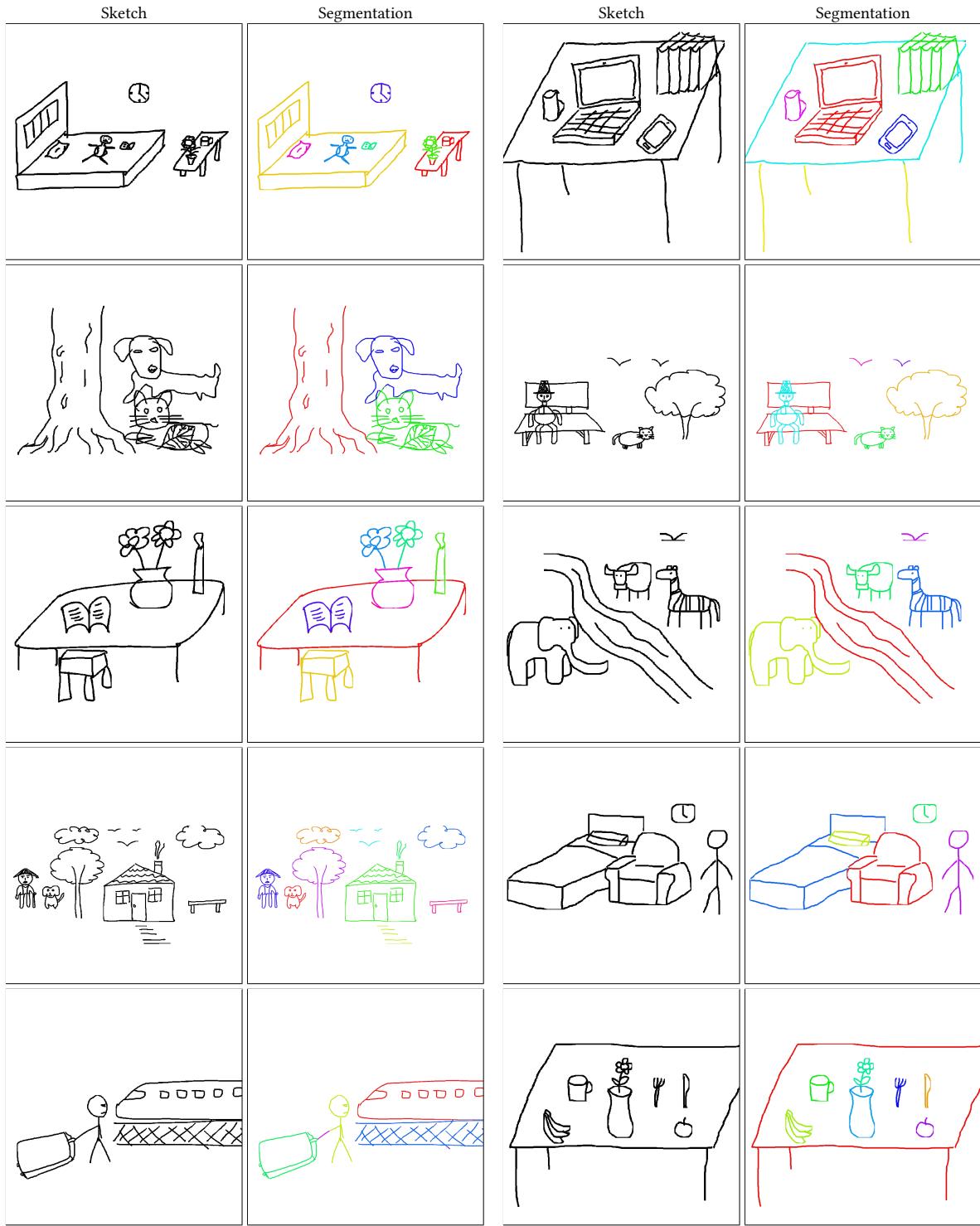


Fig. 23. **Zhang et al. Performance.** We show our segmentation outputs of Zhang et al. across a diverse range of scene sketches. The examples span indoor scenes (e.g., bedrooms, desks), outdoor urban settings (e.g., train station), and natural environments with animals, both in close-up and distant views. Our method demonstrates reasonable segmentation performance across this broad spectrum of semantic and spatial contexts.



Fig. 24. Example sketches from our InstantStyle dataset. These sketches are derived from Visual Genome [2017] containing 5 to 10 annotated objects. For visual clarity, we mask unsegmented regions in the generated sketches. This dataset contains 53 new categories that are not part of SketchyScene's 45 classes, and contain 1068 sketches in total.

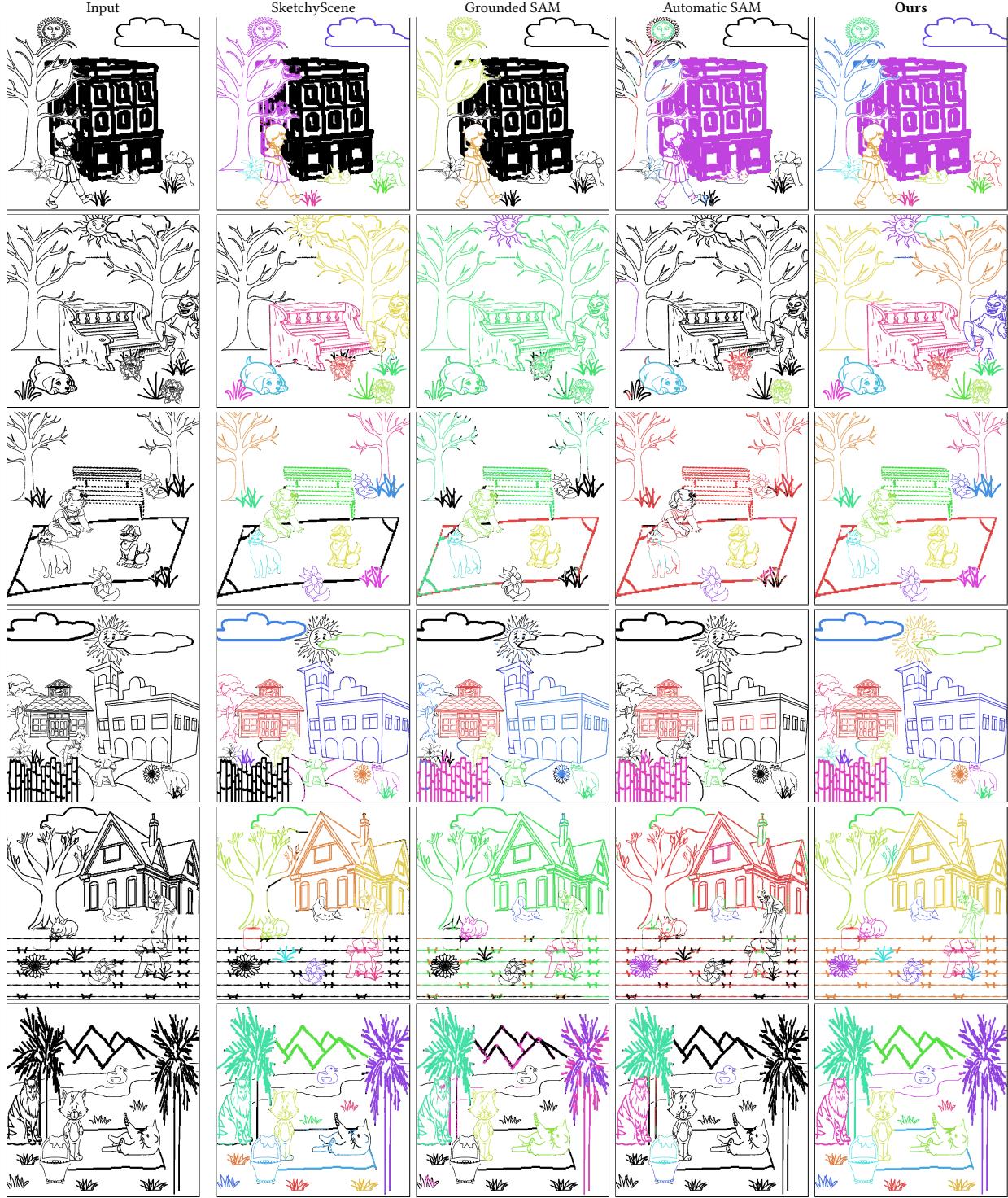


Fig. 25. Qualitative comparison of instance segmentation methods on the SketchyScene dataset. Our method surpasses SketchyScene, Grounded SAM and Automatic SAM baselines by delivering fine-grained, semantically consistent segmentations with precise boundaries. In the urban scene (row 1), our approach 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.

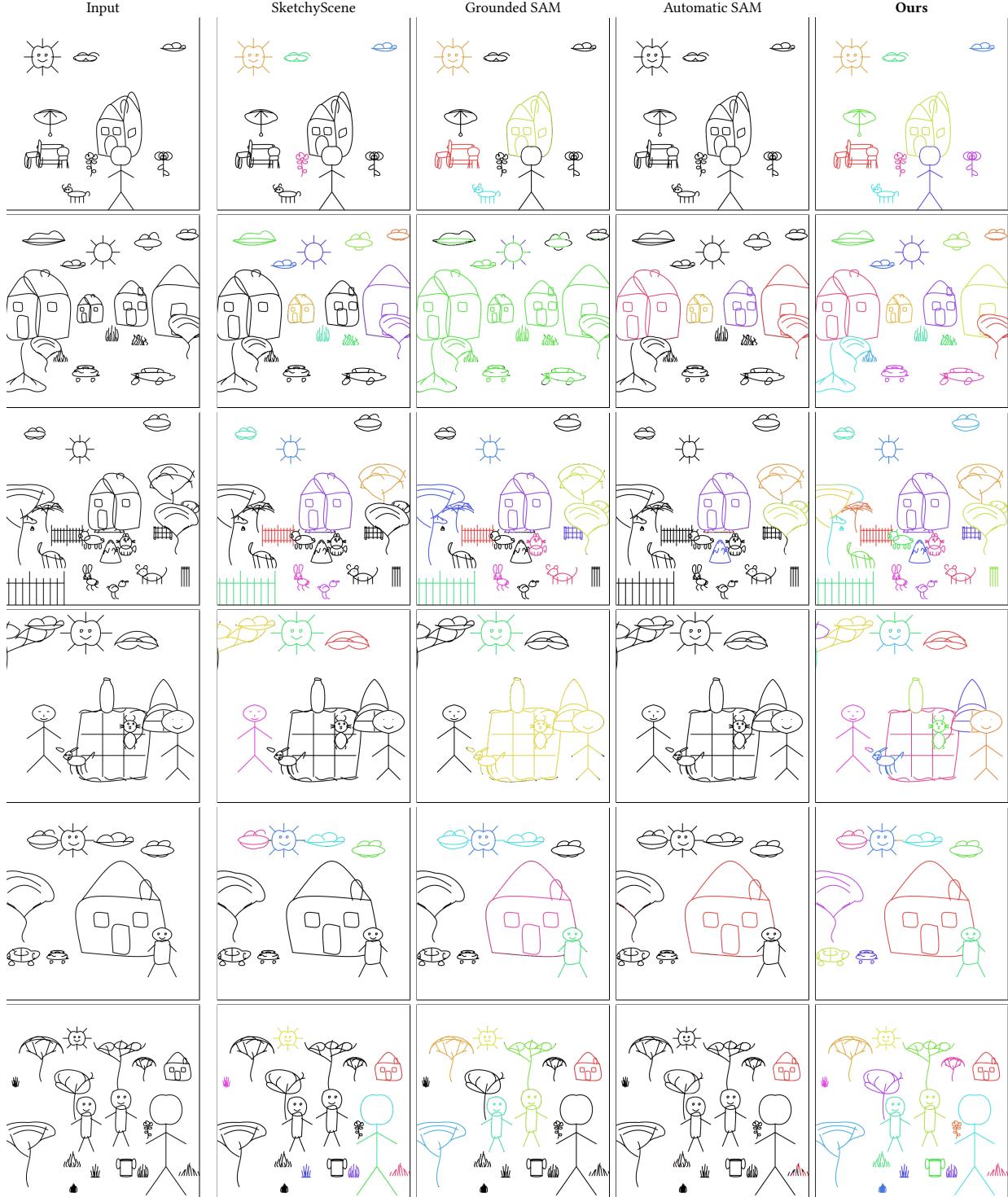


Fig. 26. Qualitative comparison of instance segmentation methods on the SketchAgent dataset. Both SketchyScene, Grounded SAM and Automatic SAM struggle to segment these abstract sketches, which differ significantly from their training data of clipart-like and real objects. Our method successfully segments individual instances while maintaining object boundaries, even in challenging cases like the last row where objects overlap with a grid-patterned picnic blanket.

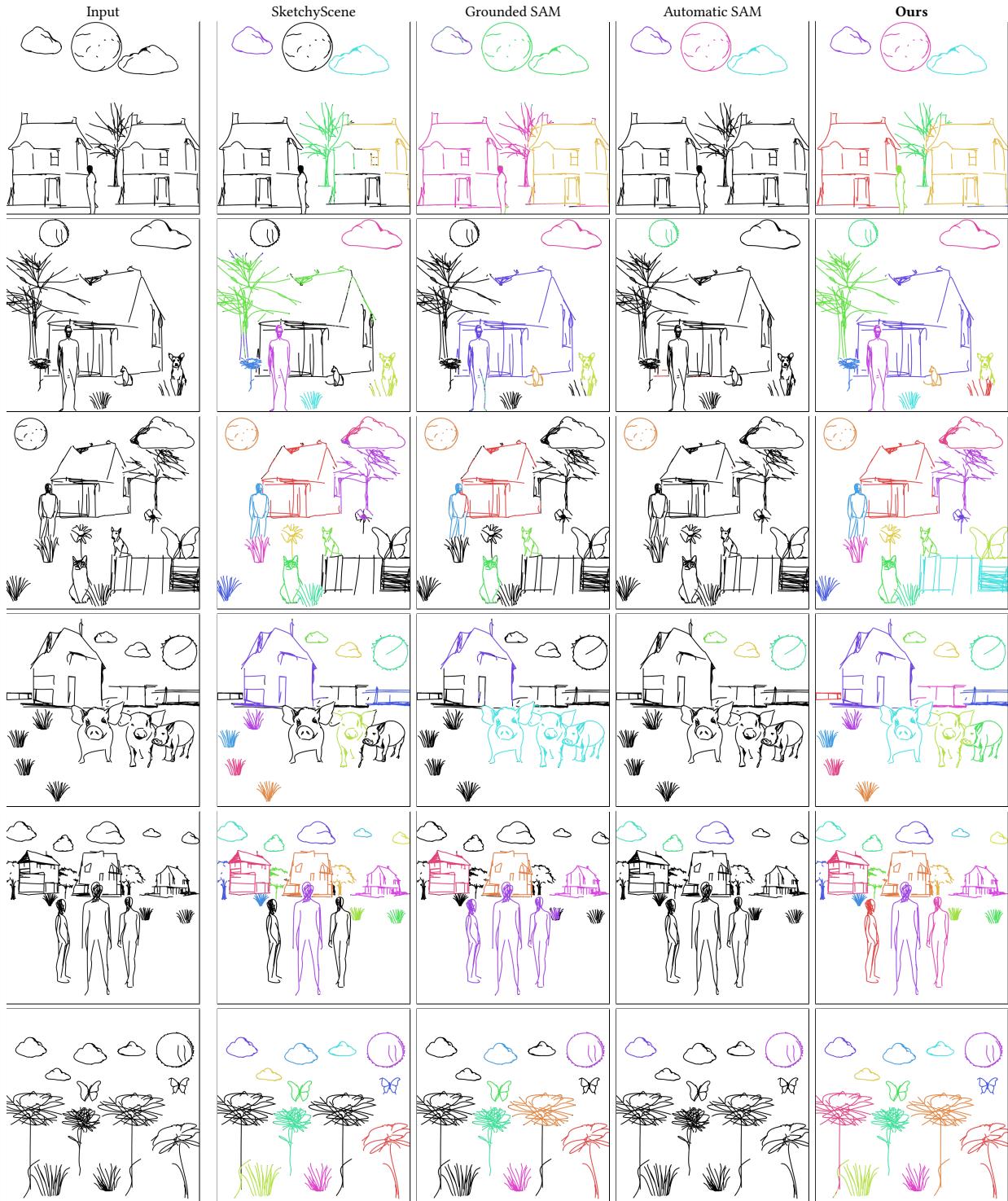


Fig. 27. Qualitative comparison of instance segmentation methods on the CLIPasso base dataset. Our method successfully detects both large and small objects with ambiguous openings in their silhouettes, whereas baseline methods either merge multiple instances into one or fail to detect them entirely.

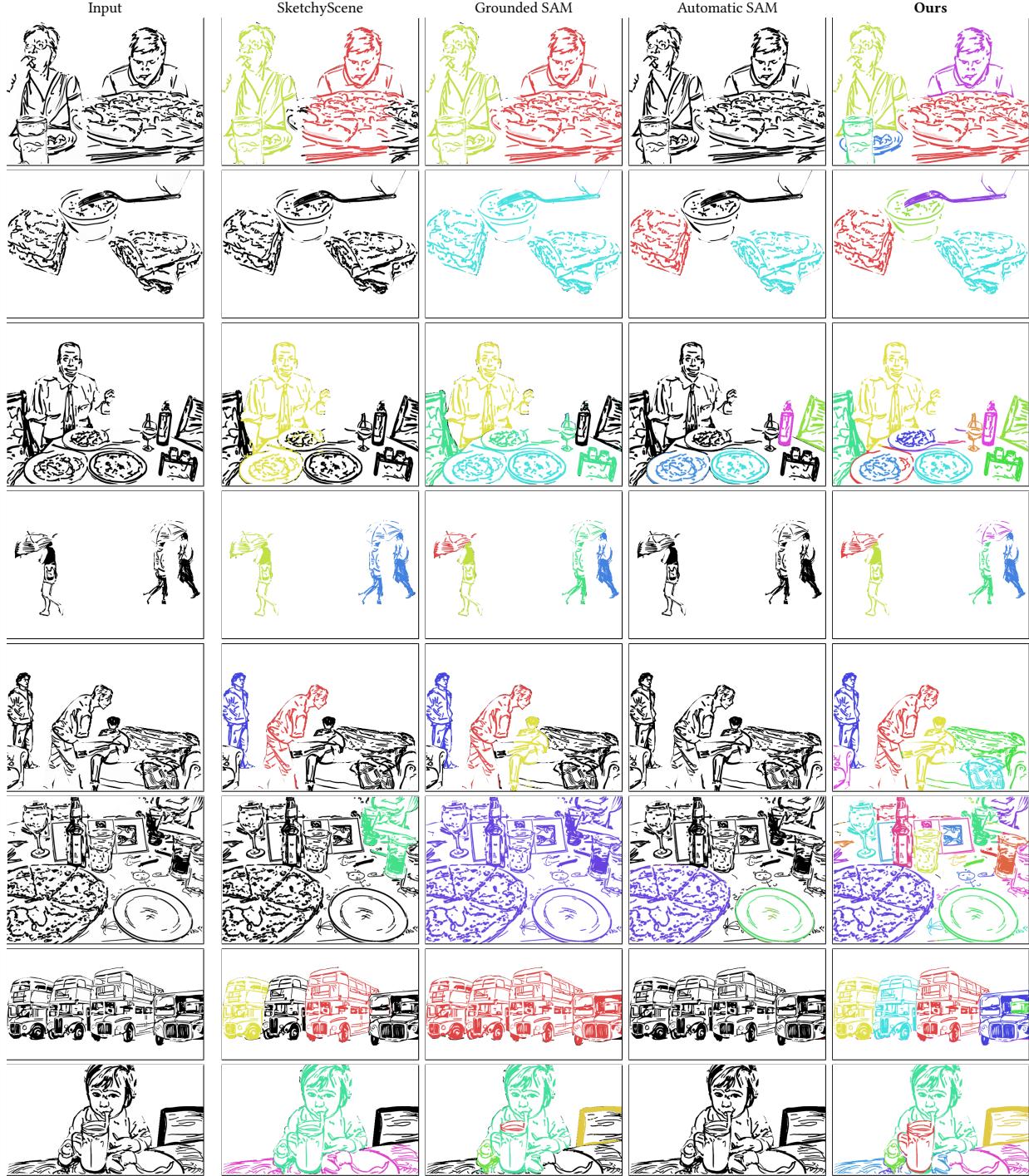


Fig. 28. Qualitative comparison of instance segmentation methods on the InstantStyle dataset. This dataset poses a significantly greater challenge compared to the others due to its increased complexity, including diverse perspectives, intricate textures, and frequent occlusions. Despite these difficulties, 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.

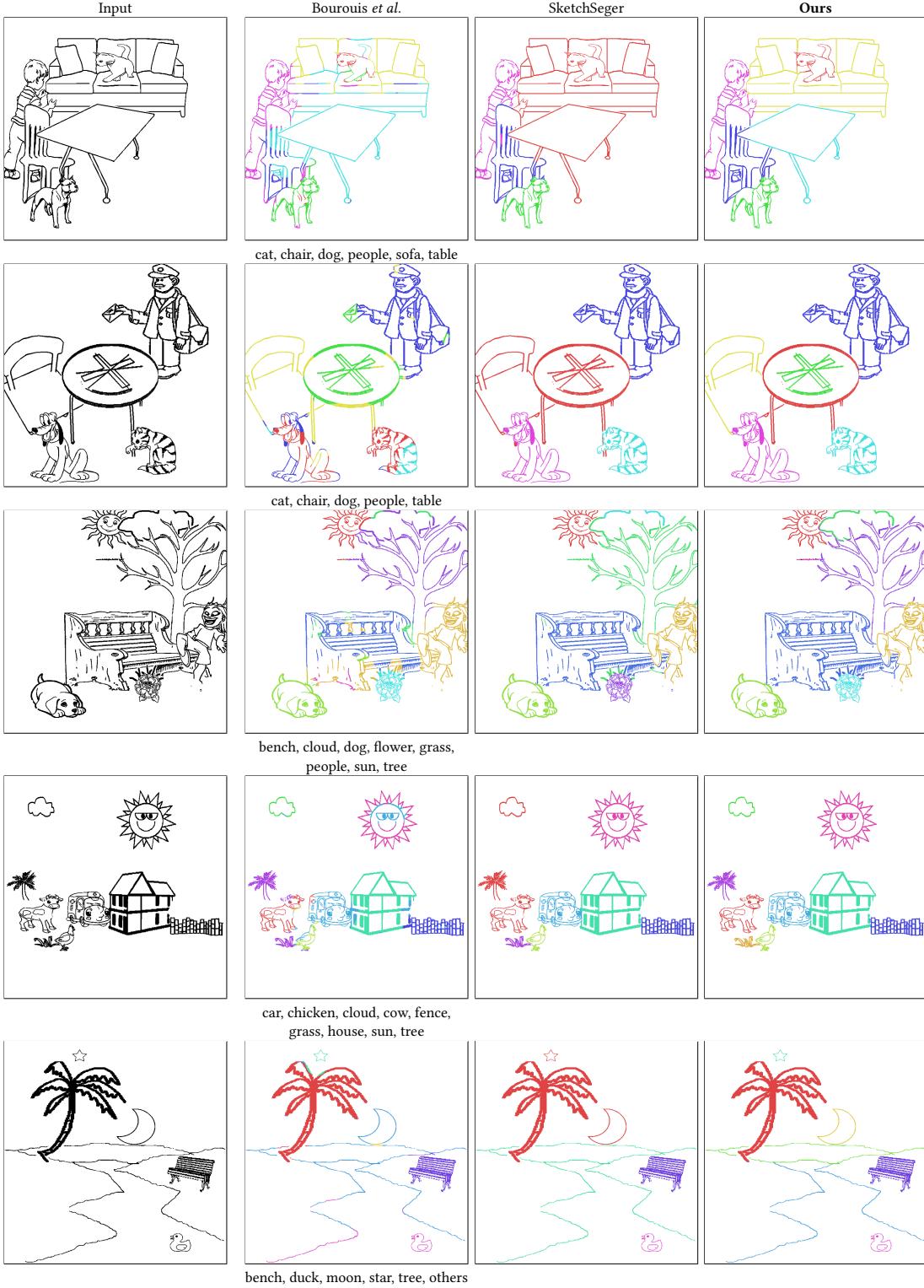


Fig. 29. **Qualitative comparison of segmentation on filtered SketchyScene dataset.** We prompt Bourouis *et al.*'s method with ground truth labels and use confidence threshold of 0.01 to ensure all sketch pixels are segmented, while using SketchSeger as-is since it does not accept input labels. However, Bourouis *et al.* struggles to accurately identify instances of the correct class, often assigning multiple class labels to the same object, as indicated by the color gradients. SketchSeger also fails to provide clean segmentations for many scenes. In contrast, our method effectively segments object instances, ensuring clear separation and consistent labeling.

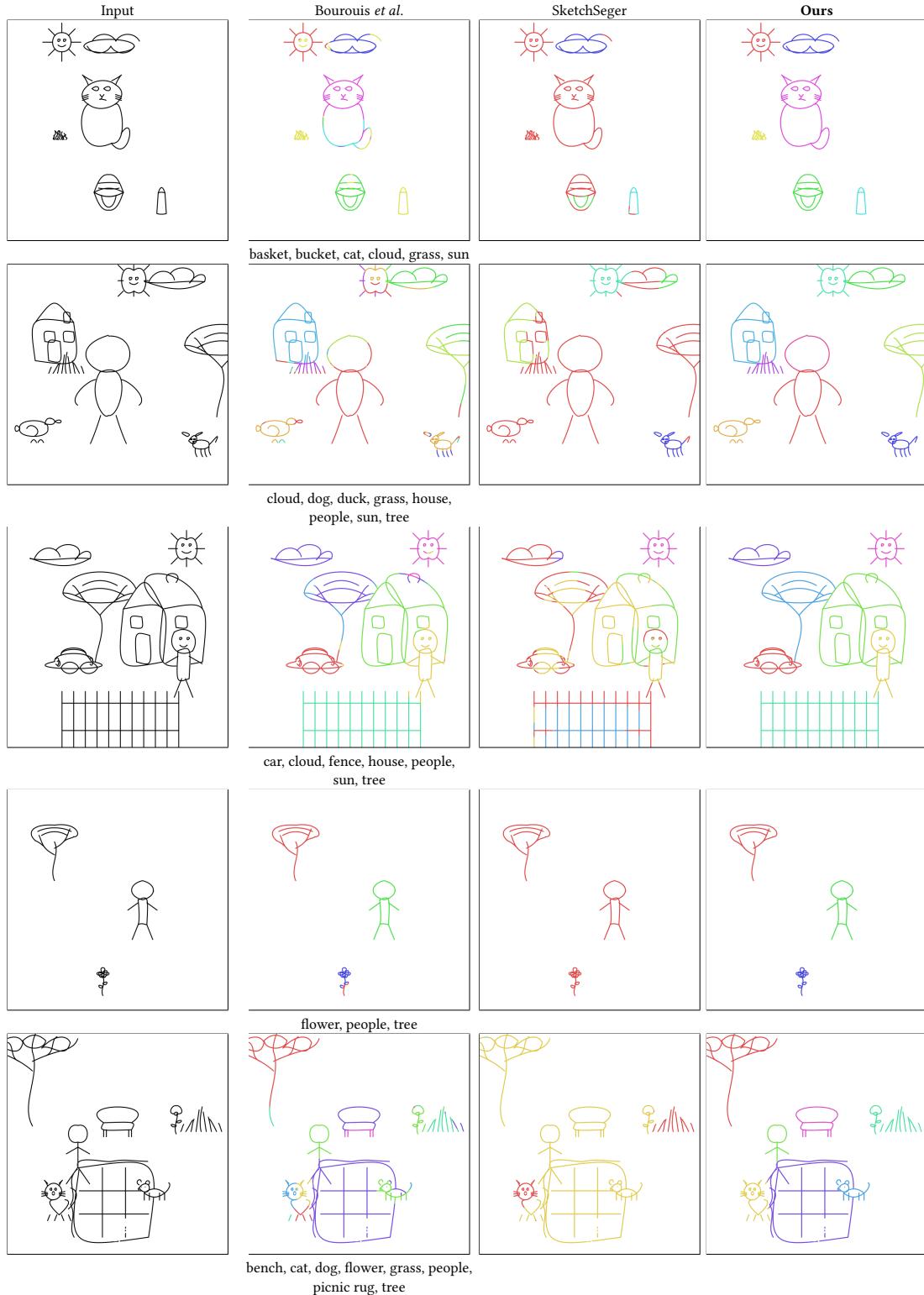


Fig. 30. **Qualitative comparison of segmentation on filtered SketchAgent dataset.** We prompt Bourouis *et al.*'s method with ground truth labels and use confidence threshold of 0.01 to ensure all sketch pixels are segmented, while using SketchSeger as-is since it does not accept input labels. However, Bourouis *et al.* struggles to accurately identify instances of the correct class, often assigning multiple class labels to the same object, as indicated by the color gradients. SketchSeger also fails to provide clean segmentations for many scenes. In contrast, our method effectively segments object instances, ensuring clear separation and consistent labeling.

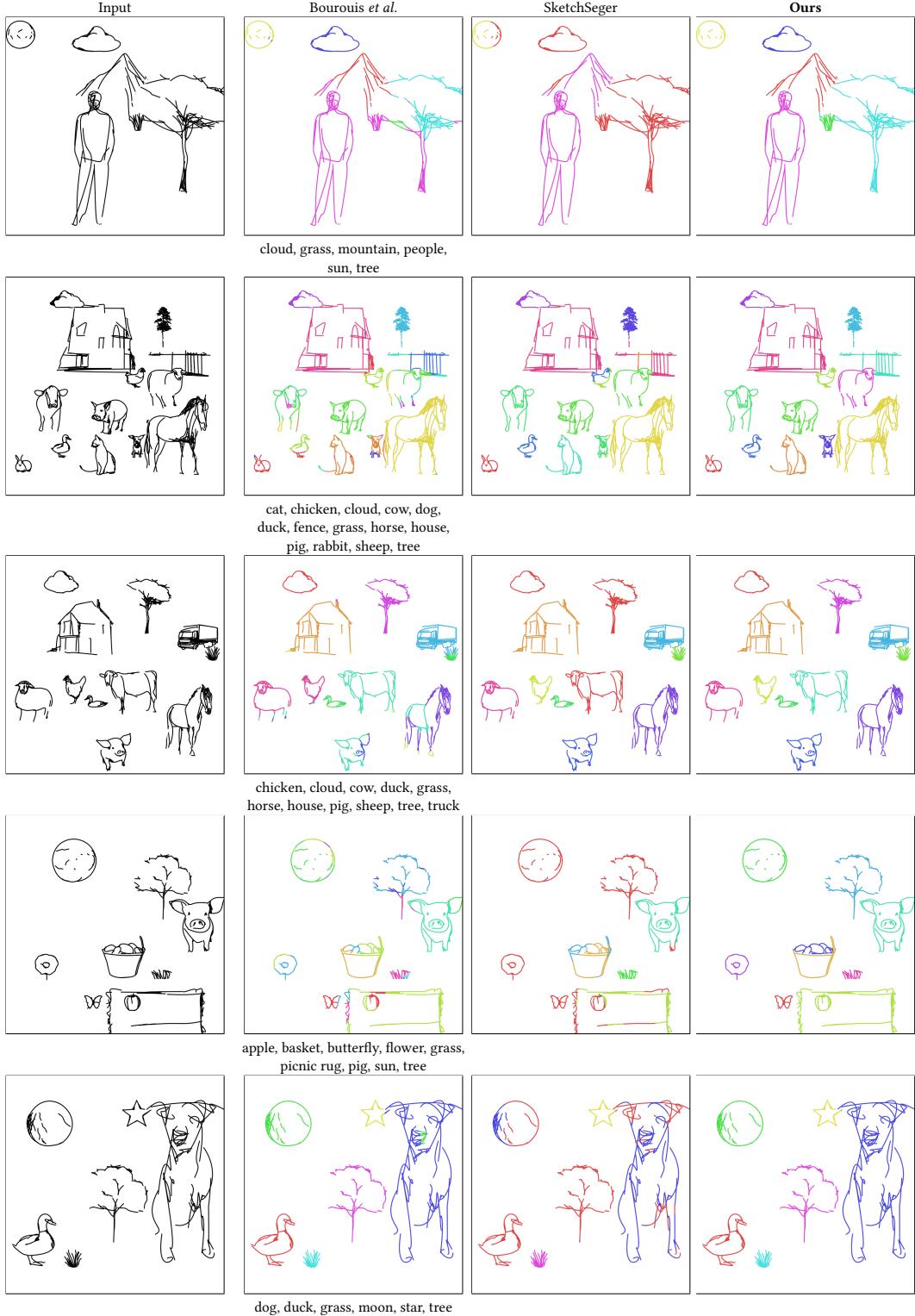


Fig. 31. Qualitative comparison of segmentation on filtered CLIPasso dataset. We prompt Bourouis *et al.*'s method with ground truth labels and use confidence threshold of 0.01 to ensure all sketch pixels are segmented, while using SketchSeger as-is since it does not accept input labels. However, Bourouis *et al.* struggles to accurately identify instances of the correct class, often assigning multiple class labels to the same object, as indicated by the color gradients. SketchSeger also fails to provide clean segmentations for many scenes. In contrast, our method effectively segments object instances, ensuring clear separation and consistent labeling.

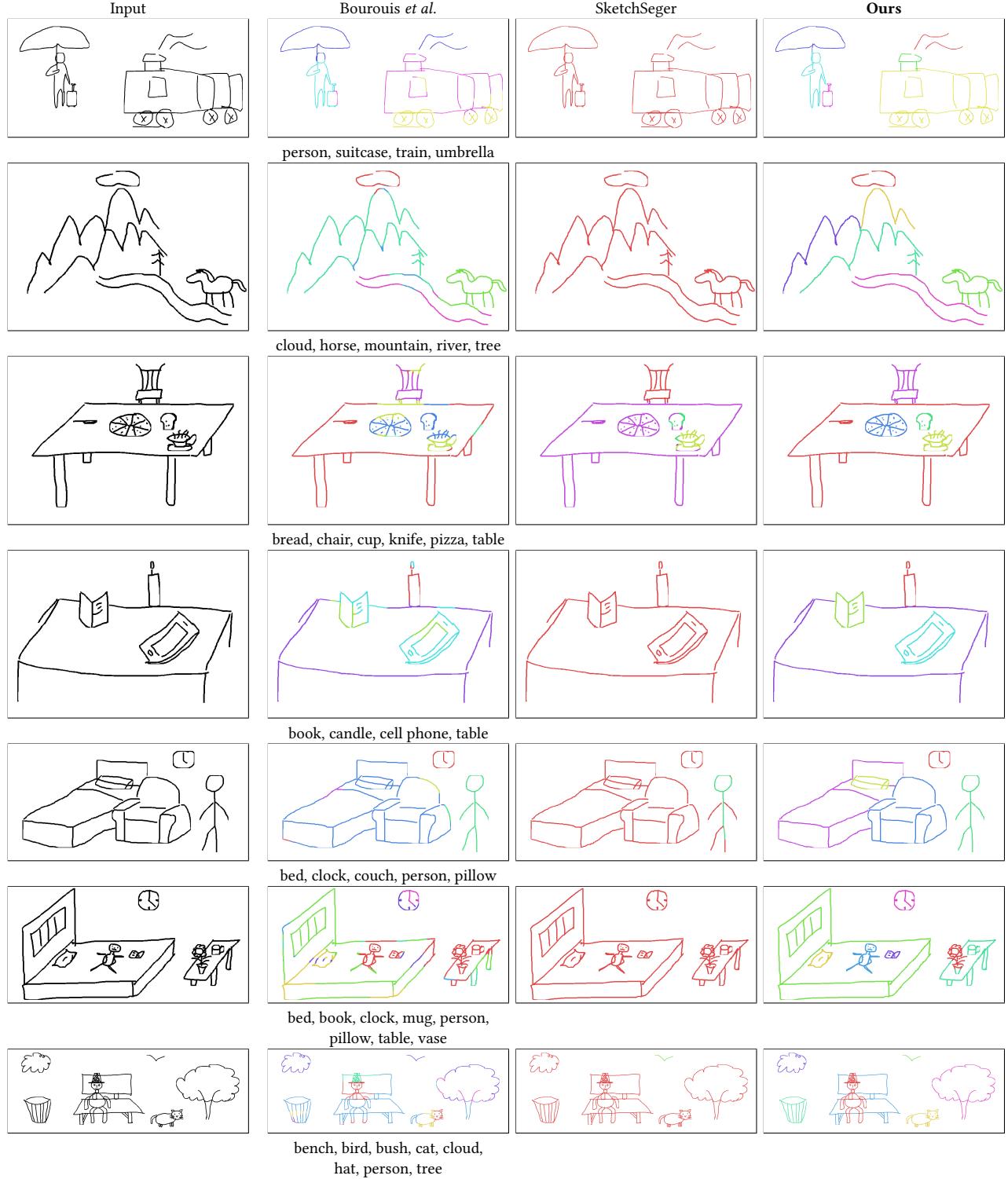


Fig. 32. **Qualitative comparison of segmentation on filtered Zhang *et al.* dataset.** We prompt Bourouis *et al.*’s method with ground truth labels and use confidence threshold of 0.01 to ensure all sketch pixels are segmented, while using SketchSeger as-is since it does not accept input labels. Our method accurately segments object instances across diverse scenes, including city life (first and last rows), natural environments (second row), and indoor settings (third to sixth rows). In contrast, baseline methods often struggle to produce clean segmentations for individual objects.

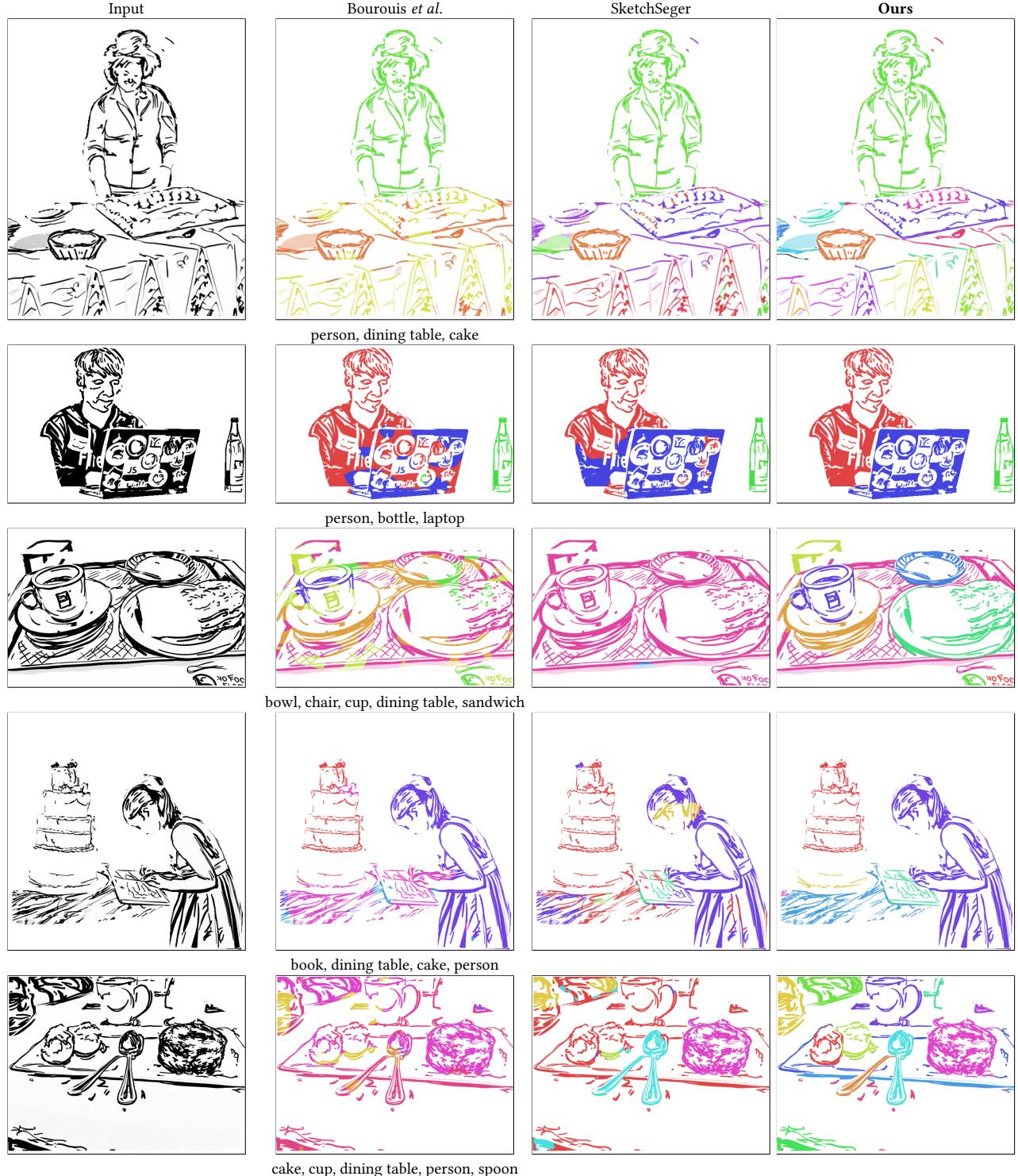


Fig. 33. **Qualitative comparison of segmentation on filtered InstantStyle dataset.** We prompt Bourouis *et al.*'s method with ground truth labels and use confidence threshold of 0.01 to ensure all sketch pixels are segmented, while using SketchSeger as-is since it does not accept input labels. Compared to Bourouis *et al.*, which sometimes fails to delineate smaller or overlapping objects, and SketchSeger, which often merges semantically distinct regions or misses foreground elements, our method produces cleaner, more precise instance boundaries and better preserves object semantics—particularly in scenes with clutter, occlusion, or fine-grained details.