

Static Segmentation by Tracking: A Label-Efficient Approach for Fine-Grained Specimen Image Segmentation

Zhenyang Feng¹, Zihe Wang¹, Jianyang Gu¹, Saul Ibaven Bueno¹, Tomasz Frelek¹, Advikaa Ramesh¹, Jingyan Bai¹, Lemeng Wang¹, Zanning Huang¹, Jinsu Yoo¹, Tai-Yu Pan¹, Arpita Chowdhury¹, Michelle Ramirez¹, Elizabeth G. Campolongo¹, Matthew J. Thompson¹, Christopher G. Lawrence², Sydne Record³, Neil Rosser⁴, Anuj Karpatne⁵, Daniel Rubenstein², Hilmar Lapp⁶, Charles V. Stewart⁷, Tanya Berger-Wolf¹, Yu Su¹, Wei-Lun Chao¹

¹The Ohio State University, ²Princeton University, ³University of Maine, ⁴University of Miami, ⁵Virginia Tech, ⁶Duke University, ⁷Rensselaer Polytechnic Institute

<https://github.com/Imageomics/SST>

Abstract

We study image segmentation in the biological domain, particularly trait segmentation from specimen images (*e.g.*, butterfly wing stripes, beetle elytra). This fine-grained task is crucial for understanding the biology of organisms, but it traditionally requires manually annotating segmentation masks for hundreds of images per species, making it highly labor-intensive. To address this challenge, we propose a label-efficient approach, **Static Segmentation by Tracking (SST)**, based on a key insight: while specimens of the same species exhibit natural variation, the traits of interest show up consistently. This motivates us to concatenate specimen images into a “pseudo-video” and reframe trait segmentation as a **tracking** problem. Specifically, SST generates masks for unlabeled images by propagating annotated or predicted masks from the “pseudo-preceding” images. Built upon recent video segmentation models, such as Segment Anything Model 2, SST achieves high-quality trait segmentation with only **one labeled image per species**, marking a breakthrough in specimen image analysis. To further enhance segmentation quality, we introduce a **cycle-consistent loss** for fine-tuning, again requiring only one labeled image. Additionally, we demonstrate the broader potential of SST, including one-shot instance segmentation in natural images and trait-based image retrieval.

1 Introduction

Understanding the sources and patterns of intra-specific variation in traits (*e.g.*, morphological characteristics such as fin length in fish or wing size in beetles) is a central goal of evolutionary and ecological study [11, 19]. Intra-specific trait variation provides a currency for assessing the roles of abiotic and biotic processes in community assembly, as it reflects the mechanisms driving species occurrence and their responses to change [76]. Museum specimens present an untapped resource for curating information on intra-specific trait variation in species morphology. Up until now, it has been difficult to harvest trait information from museum specimens due to the sheer amount of manual labor needed to make such measurements. Automatic segmentation of morphological traits from specimen images has the potential to scale up the measurement of traits and free up researchers to focus on analysis and interpretation. This paper originated from an interdisciplinary collaboration between biologists and computer scientists, aiming to segment images of organismal specimens to measure variation in traits to fill this much-needed knowledge gap.

Training a segmentation model [12, 21, 37, 41, 50] is arguably the most straightforward approach to this problem. However, it requires annotating traits on tens, if not hundreds, of images per species

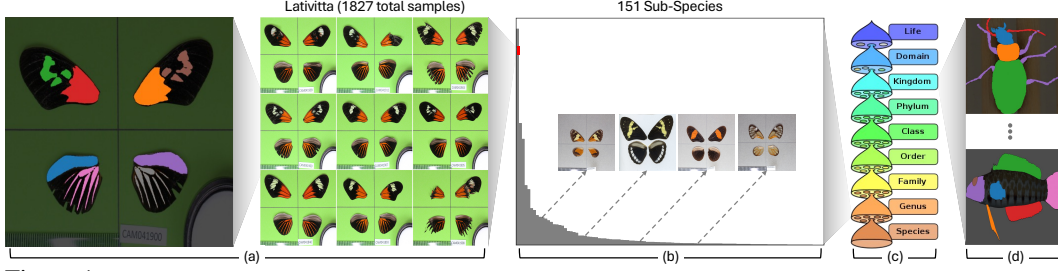


Figure 1: **Illustration of the trait segmentation problem from specimen images.** (a) Specimen samples of *Heliconius erato lativitta*, and one example of segmentation masks. (b) The histogram of sample counts per subspecies in the Cambridge Butterfly Collection [35], with exemplar images. (c) These subspecies belong to the genus *Heliconius*, which falls under the suborder *Rhopalocera*, encompassing over 10,000 butterfly species worldwide. (d) Trait segmentation is also important for other animals, such as beetles and fish.

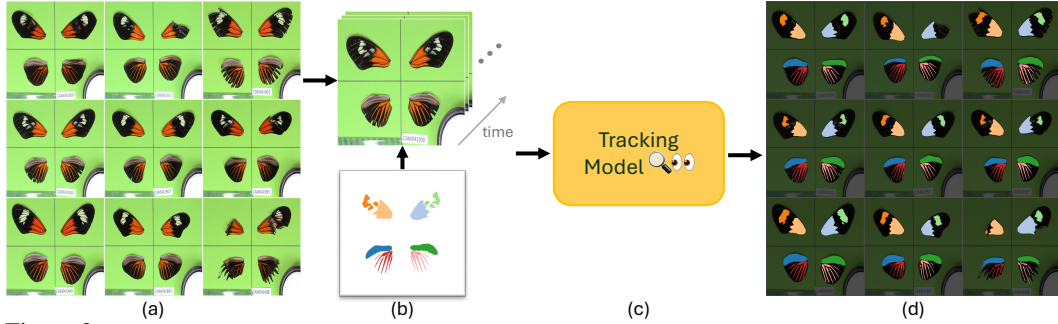


Figure 2: **Illustration of our approach Static Segmentation by Tracking (SST).** (a) Different specimens of the same species. (b) We concatenate these static, non-sequential images into a pseudo-video. (c) The annotated masks of the first image are treated as the prompt to a tracking algorithm, such as SAM 2 [69]. (d) SST can achieve high-quality trait segmentation in a one-shot manner.

to ensure the model generalizes well. This process is itself laborious, let alone there are millions of species on Earth and many of them do not have sufficient samples for labeling (see Fig. 1). Several recent segmentation algorithms focused on a few-shot setting, aiming to adapt models to new concepts using only a handful of labeled examples [25, 75, 79, 80]. However, most of these methods were designed to segment a single concept at a time (*e.g.*, an entire beetle) rather than multiple traits jointly (*e.g.*, the beetle’s head, antennae, and elytra). Even when segmenting a single trait, they often struggle to capture fine details, performing much worse than many-shot methods (see Section 4). We thus ask,

*How can we perform fine-grained segmentation on specimen images
without a large amount of labeled data?*

We begin with a deeper look at specimen images, particularly those of the same species. We make several key observations (see Fig. 1). From a *macro* perspective, where a specimen is viewed as a “whole,” biological variations naturally cause specimens of the same species to appear different; some may even have damaged parts. However, from a *micro* perspective, where a specimen is seen as a “composition” of traits—the *components we aim to segment*—specimens of the same species look quite similar in their trait layouts. Unless damaged, these traits consistently appear and maintain structured spatial relationships with one another. Notably, each trait has distinct characteristics, such as color, shape, size, pattern, and relative position, offering rich cues for identifying and locating them across specimens of the same species.

Building on these insights, we propose **reframing trait segmentation in specimen images as a tracking problem**. Tracking involves identifying an initial set of instances, assigning each a unique ID, and following them across video frames [42, 85]. *In our case, the instances are distinct traits, each marked by a unique color as in Fig. 2 (b).* While we do not have a true video, but rather a set of static specimen images, the variations observed in traits across images—such as changes in size, location, orientation, shape, and color—closely resemble the transformations seen in video frames due to camera movement, motion, deformation, and lighting changes. Even damaged parts can be viewed as occlusions in this analogy. This motivates us to concatenate *static, non-sequential* specimen

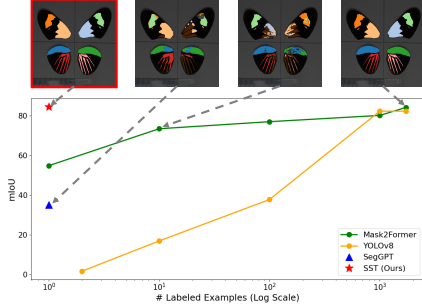


Figure 3: **Static Segmentation by Tracking (SST)** outperforms other one/many-shot baselines (on *Heliconius erato lativitta*).

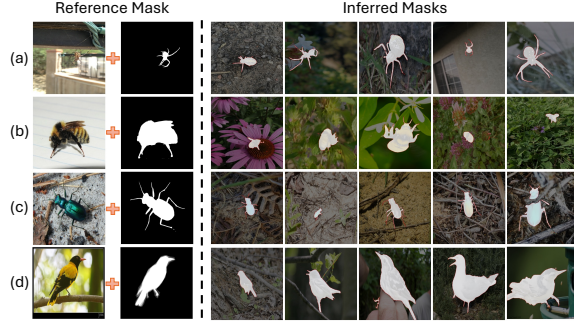


Figure 4: **SST applies to in the wild images.** (a–c) Spiders, bees, and beetles in iNaturalist [26]. (d) Birds in CUB [78].

images into a “pseudo-video” and apply a tracking algorithm to locate and segment individual traits across frames—**given only the annotated segmentation masks from the first frame** (Fig. 2).

We name our approach **Static Segmentation by Tracking (SST)**, which *lifts an image segmentation problem into a tracking problem, leveraging the latter’s characteristics to achieve the former in a remarkably label-efficient, one-shot manner*. Essentially, the model’s task is to localize each annotated mask from the first frame in subsequent frames and then *propagate* and *deform* it accordingly.

We implement SST using recent pre-trained video segmentation models, including Cutie [16], DEVA [15], and SAM 2 [69]. Given the annotated masks from the first frame as prompts, these models are capable of tracking them across frames. We evaluate SST on three specimen image datasets: Cambridge Butterfly [35], NEON Beetle [22], and Fish-Vista [47]. SST demonstrates much better performance than other one-shot baselines, such as SegGPT [80], in trait segmentation. Surprisingly, in some scenarios, SST even outperforms segmentation models trained with abundant labeled data, including Mask2Former [14] and YOLOv8 [33] (see Fig. 3). *We attribute this success to the fact that SST does not treat labeled and unlabeled images as IID samples—an assumption underlying most image segmentation algorithms—but instead explicitly leverages their dependency to facilitate segmentation.* We view this as a breakthrough in the analysis of specimen images.

Further improvement (Section 3.2). SST uses pre-trained models, and as such, it may struggle when transitions between static specimen images differ significantly from those in the training videos. To address this, we propose the **Opening-Closing Cycle-Consistent Loss (OC-CCL)** for *semi-supervised model fine-tuning*, leveraging the same labeled image as the prompt for supervision.

Further exploration (Section 4.4). We explore additional application scenarios for SST. Beyond specimen images, we also find that SST performs well on **instance segmentation of animals** in natural images, even when the “pseudo-video” contains rapidly changing and arbitrarily varying backgrounds (Fig. 4). Additionally, we investigate **trait-based image retrieval**, aiming to retrieve specimen images that share a specific trait with the query image. By repurposing OC-CCL, we demonstrate that SST can accurately retrieve images containing specific traits, such as the white bands on butterfly forewings (Fig. 9b).

Contribution. In addition to SST and OC-CCL, we present following contributions: We hand-labeled over 813 butterfly specimen images and semi-automatically labeled 2, 831 images with trait masks, covering more than 150 subspecies. We also hand-labeled 180 beetle specimen images. These labeled datasets are intended to serve as a testbed for future research in specimen image segmentation. We demonstrate the potential of SST in broader application scenarios, including instance segmentation of images taken in the wild and trait-based image retrieval.

Remark. This paper originated from an interdisciplinary collaboration aimed at addressing a real-world bottleneck in specimen image analysis. The primary challenge arises from the fine-grained nature of traits, the labor-intensive annotation process, and the vast diversity of species. Since species-specific methods are not scalable, we focus on identifying common properties that enable a more generalizable solution, ideally leveraging recent foundation models.

Our proposed approach, SST, embodies this principle by effectively leveraging dependencies across samples, resulting in a solution that is both simple and generalizable to tackle this long-standing challenge. While SST might appear straightforward *in hindsight*, its development was far from

trivial. Conventional image segmentation models remain the dominant approach, yet they struggle with scalability and adaptability across species. The key novelty of our work lies in recognizing and implementing the *appropriate* way to address specimen image segmentation—one that transcends the label-intensive, species-specific constraints in favor of a more generalizable framework. Please also refer to Section 5 for a discussion on the paper’s scope.

2 Related Work and Background

Image segmentation is a long-standing challenge in computer vision [12, 21, 37, 41, 50]. Semantic and instance segmentation are among the most popular tasks today, aiming to group pixels with the same semantic meanings and instances [23, 50]. While much of the focus has been on segmenting common, coarse-grained objects, recent works have begun exploring part-level segmentation within these objects [28, 66]. In this paper, we address an underexplored challenge: trait segmentation for fine-grained object categories, such as subspecies of butterflies.

While state-of-the-art (SOTA) models have shown impressive capabilities in segmentation [14, 33], collecting pixel-level annotations for training is labor-intensive. To address this, **few-shot segmentation (FSS)** has emerged as a promising paradigm, attempting to use few-shot learning techniques to generate high-quality segmentation masks for new classes [24, 40, 67, 75, 86]. Here, we propose a novel perspective and algorithm to tackle the FSS problem.

Video segmentation focuses on segmenting the same concepts (*e.g.*, objects) across video frames [84]. Compared to image segmentation, it requires associating masks between frames to assign the same labels. Recent models [15, 16, 69] are mainly built upon the transformer architecture with memories, and trained on a large-scale video data. Given annotated masks from the first frame, they can track and segment the target instances in subsequent frames.

While image and video segmentation have typically been studied separately, we show that models developed for the latter can be effectively applied to the former in a label-efficient manner, even when the images are non-sequential.

Co-segmentation and image registration share similar properties with our approach—they leverage dependencies across images. Co-segmentation locates objects that appear in multiple images in an unsupervised manner [13, 17, 34, 70, 77]. Our work can be viewed as a one-shot supervised approach, leveraging tracking models to efficiently segment the traits of interest across images. Image registration is widely used to densely align pixels between images, such as brain MRI scans [49]. In our application, we do not need exact pixel-to-pixel associations; we only need to localize the traits of interest across images.

3 Proposed Approach

Problem definition and notation. We study trait segmentation from specimen images of the *same* species. Let $\mathbf{x} \in \mathbb{R}^{W \times H \times 3}$ denote a $W \times H$ image and $\mathbf{y} \in \{0, 1\}^{W \times H \times C}$ denote the corresponding ground-truth segmentation masks of C distinct traits. The goal is to develop a segmentation model f such that its output $\hat{\mathbf{y}} = f(\mathbf{x})$ matches \mathbf{y} .

Typically, one needs to collect a labeled training set with ample pairs of (\mathbf{x}, \mathbf{y}) , and use it to train f in a supervised way. In this paper, we target the one-shot scenario, *i.e.*, building f using a single labeled image (\mathbf{x}, \mathbf{y}) .

3.1 Static Segmentation by Tracking (SST)

At first glance, this seems like a daunting challenge. However, the domain-specific properties described in Section 1 provide a crucial foundation. Our proposed approach, **SST**, leverages these properties by reframing trait segmentation as a tracking problem, which naturally becomes a one-shot task given a set of labeled instances in the first frame.

More specifically, let $\{\mathbf{x}_0, \dots, \mathbf{x}_N\}$ denote a *sequence* of video frames; a tracking algorithm aims to track the instances in \mathbf{x}_0 , encoded by the label \mathbf{y}_0 , across the remaining frames. In our context, we do not have a real video but rather a *set* of N unlabeled images and one labeled image (\mathbf{x}, \mathbf{y}) .

Nevertheless, the domain-specific properties motivate us to construct a “pseudo-video” by treating x as the first frame x_0 , followed by the unlabeled images.

SST with pre-trained models. We leverage recent pre-trained video segmentation models for tracking [15, 16, 69]. Despite differences in technical details, these models share a similar architecture. *Without loss of generality, we focus on the Segment Anything Model 2 (SAM 2) [69] for the remainder of this section.* Below, we first briefly introduce its model architecture and inference mechanism.

SAM 2 uses a promptable Transformer encoder-decoder f augmented with a memory bank B to process a video and generate masks (see Fig. 5). Let $\{y_0, \dots, y_N\}$ denote the ground-truth labels for video frames $\{x_0, \dots, x_N\}$. When the label of x_n is unavailable, we set $y_n = \emptyset$. Let \hat{y}_n denote the predicted mask for x_n , and let B_n represent the updated memory bank after the prediction, which stores both the feature and mask information. B_n can then be accessed by the next frame x_{n+1} to connect consecutive frames. In the context of tracking, B_n can be interpreted as the updated state estimate after perceiving the measurement x_n .

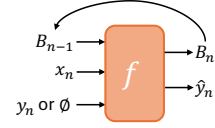


Figure 5: SAM 2 inference mechanism.

At each timestamp n , f takes the tuple $[x_n, B_{n-1}, y_n]$ as input, where y_n is treated as the (optional) prompt. It then outputs the tuple $[B_n, \hat{y}_n]$, where B_n will be used as input at the next timestamp,

$$[B_n, \hat{y}_n] = f([x_n, B_{n-1}, y_n]).$$

In our context, we have $y_n = \emptyset, \forall n > 0$, meaning only the first frame x_0 is annotated. By inputting y_0 to f at timestamp 0, we instruct the model on what to *segment*—the distinct traits and their extents. The resulting memory bank B_0 then carries this information to successive frames, guiding the model on what to *track* across frames to generate the masks $\{\hat{y}_1, \dots, \hat{y}_N\}$. See Fig. 2 for an illustration.

Pseudo-video creation. There are multiple ways to concatenate unlabeled images into a pseudo-video. Intuitively, an order with smooth transitions is preferred, as it could potentially improve SST’s performance. In contrast, non-smooth transitions may degrade SST. However, searching for an optimal order incurs additional computational cost.

To eliminate the uncertainty in creating pseudo-videos, we implement SST by constructing multiple short, two-frame videos, unless stated otherwise. Specifically, given the labeled image-mask pair (x_0, y_0) and N unlabeled images $\{x_1, \dots, x_N\}$, we construct $\{x_0, x_1\}, \dots, \{x_0, x_N\}$ and apply SST independently to each.

We note that the above implementation also ensures a fair comparison to the baselines. Conventional image segmentation models process each test sample independently, reflecting the real-world online use case where a newly captured image is processed immediately. Recent few-shot approaches similarly process each test image independently. Noting that considering all test images jointly would transform the conventional *inductive* setting into a *transductive* one, we choose to process each test image independently.

That said, in the Appendix, we explore concatenating multiple unlabeled images. In the long run, developing an approach to search for the optimal order would be valuable.

Remark. According to the original paper [69], SAM 2 is readily applicable to a batch of static, non-sequential images, by *setting the memory bank B_n to empty*. In essence, without the memory bank, SAM 2 treats each input image as an IID sample and processes them independently.

Our insight is that even if the input images are taken in a non-sequential manner, whenever there exists a useful dependency among them (*e.g.*, from the same species), SAM 2 has the potential to leverage this dependency. The key is to allow the memory bank to update, rather than resetting it.

3.2 One-Shot Fine-Tuning for SST

SST uses pre-trained video models in a plug-and-play fashion without altering their weights, even though our use case might be outside the training data distribution. As a result, SST is expected to fail when transitions between static images are significantly out-of-distribution (OOD).

One intuitive way to address this is model fine-tuning. However, with only one labeled image (x_0, y_0) , fine-tuning risks overfitting. Additionally, since we have used y_0 to prompt the model, we face another challenge: *it is unclear how to use it “dually” as the label to supervise fine-tuning.*

Opening-Closing Cycle-Consistent Loss (OC-CCL). To overcome these challenges, we propose a novel fine-tuning approach that leverages the flexibility of creating pseudo-videos. We can duplicate static images and inject them into the video sequence at different timestamps, allowing us to obtain multiple predictions for the same image. Specifically, given a short pseudo-video $\{x_0, x_1\}$, we duplicate both images, denoted by \dagger , and create a palindrome-style cycle $\{x_0, x_1, x_1^\dagger, x_0^\dagger\}$, inspired by [29]. In this cycle, the labeled image x_0 serves as both the “opening” and “closing” frames. Note that we do not require the label of x_1 .

Unlike timestamp 0, where y_0 serves as the prompt for the tracking model f , at the last timestamp, we treat x_0^\dagger as an unlabeled frame without prompts (*i.e.*, $y_0^\dagger = \emptyset$). This design allows us to use y_0 , the ground-truth label of x_0^\dagger , to supervise the fine-tuning of f —by minimizing the discrepancy between the predicted \hat{y}_0^\dagger and y_0 . The rationale is that if f fails to track traits in the intermediate frames (*i.e.*, x_1 and x_1^\dagger), it will not carry useful information for correctly segmenting x_0^\dagger (see Fig. 6).

We use a combination of binary cross entropy (BCE) loss and Dice loss for fine-tuning, both of which are commonly used in training segmentation models.

Implementation details. We assume access to one labeled image (x_0, y_0) and a set of unlabeled training images disjoint from the test images. At each fine-tuning step, we sample x_1 from the unlabeled set and create a short palindrome-style cycle. We apply LoRA [27] to fine-tune the decoder and memory encoder of video segmentation models.

Since short palindrome-style cycles $\{x_0, x_1, x_1^\dagger, x_0^\dagger\}$ are used, the memory bank can retain the prompt y_0 until the closing frame, potentially making fine-tuning ineffective. To address this, we apply the following strategy.

1. During the forward pass, we reset the memory bank after processing x_1 , preventing it from carrying y_0 to x_0^\dagger .
2. To propagate tracking information from x_1 to x_1^\dagger , we use the predicted mask \hat{y}_1 from the former as the prompt for the latter. Notably, this prompt remains differentiable.

With this strategy, minimizing OC-CCL encourages f to predict an accurate \hat{y}_1 , ensuring that it propagates useful information for correctly segmenting x_0^\dagger at the final timestamp to match the ground-truth y_0 .

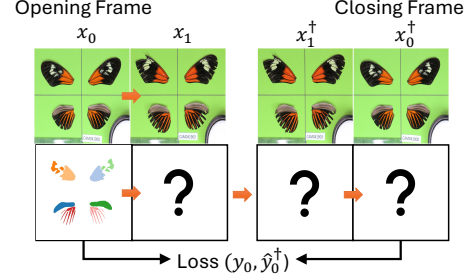


Figure 6: **Opening-Closing Cycle-Consistent Loss (OC-CCL)** compares predicted masks to the ground truth.

3.3 Extension to Trait-Based Retrieval

Beyond trait segmentation within the same species, SST can also be used to retrieve specimens exhibiting similar traits (*e.g.*, the white band on the forewing or the orange tiger tails on the hindwing) across different species. Given a query image x_0 and a target trait y_0^\star —a single channel in the original $y_0 \in \{0, 1\}^{W \times H \times C}$ —SST scores each image x_i in the retrieval pool by

1. creating a palindrome-style cycle $\{x_0, x_i, x_i^\dagger, x_0^\dagger\}$;
2. using y_0^\star as the prompt and taking the forward pass introduced in Section 3.2 to predict \hat{y}_0^\dagger ;
3. calculating the IoU between y_0^\star and \hat{y}_0^\dagger .

Namely, if x_i has the target trait, the trait mask y_0^\star should accurately propagate to x_i and then propagate back to x_0 .

Table 1: **Specimen segmentation results (mIoU) and computational cost.** SST outperforms recent FSS methods as well as standard many-shot segmentation models trained on full data, while with much less inference computational cost.

# of Data	Model	Time (s)	Major [35]	Minor [35]	Fish [47]	Beetle [22]
One-Shot	HDMNet [67]	0.53 \pm 0.00	4.2 \pm 0.8	4.0	2.4	6.1 \pm 2.0
	PFENet [75]	0.43 \pm 0.05	8.0 \pm 3.4	4.2	3.1	19.1 \pm 3.4
	VAT [24]	0.39 \pm 0.03	13.5 \pm 3.7	15.1	24.6	26.0 \pm 5.7
	SegGPT [80]	0.27 \pm 0.04	35.2 \pm 2.9	41.9	54.9	45.2 \pm 4.3
One-Shot	DEVA [15] + SST	0.13 \pm 0.04	73.1 \pm 4.0	68.6	50.8	39.0 \pm 6.9
	Cutie [16] + SST	0.11 \pm 0.04	67.4 \pm 5.0	69.7	51.9	45.8 \pm 4.4
	SAM 2 [69] + SST	0.30 \pm 0.06	81.0 \pm 1.0	70.6	70.4	61.9 \pm 3.7
Full	YOLOv8 [33]	0.51 \pm 0.12	71.1	-	-	75.1
	Mask2Former [14]	0.56 \pm 0.07	79.3	-	-	83.2

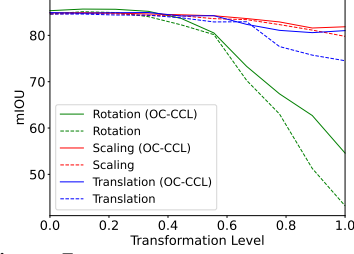


Figure 7: **Fine-tuning results.** OC-CCL fine-tuning improves SST’s robustness to transformation variations.

Table 2: **Opening-Closing Cycle-Consistent Loss and multi-shot results.** Applying OC-CCL and Multi-Shot inference on top of SST contributes to further performance improvement, surpassing SegGPT by a large margin.

Modules		SegGPT [80]		DEVA [15] + SST		Cutie [16] + SST		SAM 2 [69] + SST	
OC-CCL	Five-Shot	Major [35]	Beetle [22]	Major [35]	Beetle [22]	Major [35]	Beetle [22]	Major [35]	Beetle [22]
✓		35.2 \pm 2.9	45.2 \pm 4.3	73.1 \pm 4.0	39.0 \pm 6.9	67.4 \pm 5.0	45.8 \pm 4.4	81.0 \pm 1.0	61.9 \pm 3.7
		-	-	74.2 \pm 2.5	48.4 \pm 6.7	77.7 \pm 2.8	48.7 \pm 2.2	81.2 \pm 1.0	65.2 \pm 3.3
✓	✓	42.3 \pm 0.3	70.3 \pm 2.7	83.2 \pm 1.9	65.5 \pm 2.4	83.4 \pm 1.7	67.2 \pm 2.0	83.4 \pm 0.4	74.2 \pm 2.9
	✓	-	-	83.4 \pm 0.7	66.9 \pm 1.3	84.7 \pm 1.3	67.6 \pm 1.1	83.9 \pm 0.2	75.6 \pm 1.5

4 Experiment

4.1 Experimental Setup

Data. We evaluate SST on three specimen data sources.

- **Butterfly:** We use the Cambridge Heliconius Collection [35]¹, annotated in consultation with biologists and the field guide [1]. Due to its long-tailed distribution (see Fig. 1), we divide the dataset into Major and Minor parts. The Major part comprises the five largest subspecies, with 100 specimens per subspecies for testing and the remaining 2, 831 samples for training on a subspecies basis. The Minor part consists of 146 subspecies with 2 \sim 3 hand-labeled samples for each of them. The training set cannot be constructed for the Minor part due to insufficient samples per subspecies.
- **Fish:** We use the Fish-Vista dataset [47], containing specimens from over 1, 900 species.² A subset (1, 573 images across 474 species) was labeled with 9 expert-selected body parts; the labels are consistent across species. Similar to Butterfly Minor, there is no training set for Fish.
- **Beetle:** We use the individual image subset of the 2018 NEON-beetles dataset [22]. We hand-labeled 180 specimen images (120/60 for training/testing) with 5 body parts. The data is challenging due to up to 90-degree of body rotations, missing, and overlapped parts.

We note that our task imposes challenges from its fine-grained nature and the vast diversity of species. As such, standard few-shot datasets may not be ideal for development. Additional dataset details are provided in the Appendix.

Evaluation metric. We use mean IoU (mIoU) as the main metric, averaged over traits in images. For FSS methods (including SST) in a one-shot setting, one labeled image (with canonical shapes and visually clear traits) is sampled as the reference for the test set. We report averaged mIoU over 20 runs with standard deviation, unless stated otherwise.

Implementation details of SST. We apply SST to three video segmentation models, DEVA [15], Cutie [16], and SAM 2 [69], using the official pre-trained models.

Baselines. We consider four representative few-shot segmentation (FSS) methods, PFENet [75], VAT [24], HDMNet [67], and SegGPT [80]. We also apply two representative many-shot instance segmentation algorithms, YOLOv8 [33] and Mask2Former [14], whenever we have sufficient training samples. The Appendix presents more baseline settings and comparisons.

Table 3: **Instance segmentation results on CUB.** Our method SST achieves similar results as SOTA FSS methods on object instance segmentation.

Model	One-Shot	Five-Shot
HDMNet [67]	65.8 \pm 1.4	66.3 \pm 1.4
PFENet [75]	72.2 \pm 0.6	73.1 \pm 0.4
VAT [24]	83.4 \pm 1.0	85.3 \pm 0.7
SegGPT [80]	51.3 \pm 2.9	78.8 \pm 1.4
SAM 2 + SST (ours)	71.1 \pm 1.3	77.9 \pm 0.8
SAM 2 + SST + OC-CCL (ours)	77.8 \pm 1.1	79.6 \pm 0.4

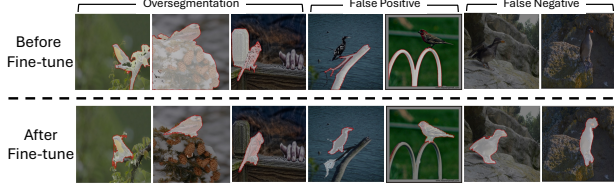


Figure 8: **Before vs. after OC-CCL fine-tuning.** Fine-tuning with OC-CCL notably improves SST with merely one labeled training example.

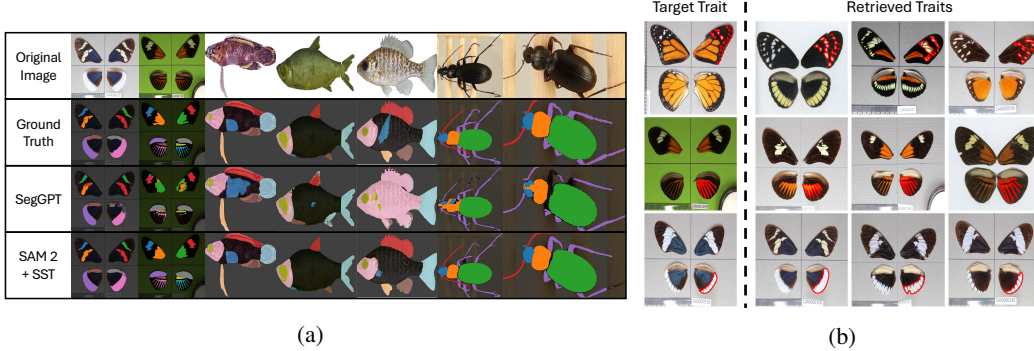


Figure 9: (a) **Qualitative results:** Trait segmentation of SST (with SAM 2) vs. SegGPT on butterfly, fish, and beetle data. (b) **Trait-based retrieval.** We can retrieve different butterfly subspecies with visually consistent traits to the specified target (marked in red).

4.2 Main Result

Specimen segmentation results. We evaluate SST on Butterfly [35] (Major, Minor), Fish [47], and Beetle [22] datasets. On Butterfly Major and Beetle with sufficient training data, we train Mask2Former [14] and YOLOv8 [33] models using all the training samples per (sub)species as many-shot baselines. For FSS algorithms and SST, we consider a one-shot setting. On the remaining datasets (Butterfly Minor, Fish), where samples per (sub)species are limited, we evaluate only in the one-shot setting in a leave-one-test-sample-out manner; standard deviation is not reported.

As shown in Table 1, SST outperforms existing FSS algorithms across most datasets using various video segmentation models. Especially on Butterfly Minor, SST achieves a margin of at least 26.7 mIoU over SegGPT [80]. Remarkably, SST based on SAM 2 [69] even surpasses many-shot algorithms trained with at least 150 samples per subspecies on Butterfly Major. As evidenced in Fig. 9a, SST offers more accurate segmentation results than SegGPT [80].

The superior performance of SST can be primarily attributed to its alignment with the task. Fine-grained specimen segmentation is inherently more challenging than the tasks for which existing FSS models were designed. However, specimen images exhibit strong interdependencies, even when captured non-sequentially, making video segmentation particularly well-suited for this problem.

Opening-Closing Cycle-Consistent Loss. Given sufficient data from Butterfly Major [35] and Beetle [22], we further fine-tune the pre-trained models with OC-CCL. For each species, one labeled image and the other unlabeled images from the training set are used to construct the palindrome cycles. We evaluate the effectiveness of OC-CCL in Table 2, where it is applied to all three tracking models (the second row), with consistent improvement across all tasks.

Multi-shot inference. We have investigated the one-shot case where only one labeled sample is leveraged in Table 1. Given more labeled samples, SST also supports multi-shot inference with more comprehensive information. As shown in Table 2, using multiple reference frames enables SST to outperform YOLOv8 [33] and Mask2Former [14], while boosting Cutie [16] to performance levels comparable with SAM 2 [69] on Butterfly Major. Moreover, combining multi-shot inference

¹Sources: [30–32, 45, 46, 48, 51–65, 68, 71–74, 81–83].

²This dataset is comprised of specimen images from various collections: [2–10, 18, 20, 43, 44].

with OC-CCL yields considerable performance improvement across all metrics. Please refer to the Appendix for more details of the multi-shot setting.

4.3 Analysis

Out-of-distribution (OOD) robustness. The actual application of SST on specimen images might encounter OOD cases, where the specimens are not captured in standard views. That is, the images might be subjected to rotation, translation, or scaling. Accordingly, we manually apply these transformations to the Butterfly test set to create OOD cases. We define transformation levels from 0.0 to 1.0, corresponding to no transformation and the largest degree of transformation we apply, respectively. At level 1.0, we randomly rotate the image between -90° and 90° , translate it up to 60% of its height or width, and scale it down to 50% of its original size. We found that SST tends to lose track of the fine-grained details after a certain level of transformations, likely due to the absence of such huge variations (between consecutive frames) in the pre-training data.

To address this, we use OC-CCL to fine-tune the model in a one-shot setting for each test image. OC-CCL consistently improves SST’s robustness as seen in Fig. 7. In the extreme rotation cases (the right end of the figure), we boost the mIoU from 40% to over 50%, a more than 10% gain.

4.4 Extension and Further Exploration

Instance segmentation. Besides fine-grained specimen segmentation, SST can also be applied to standard object instance segmentation on images taken in the wild. We use the CUB-200-2011 dataset [78] to demonstrate such a capability. Given one/five random bird images and their segmentation masks, we examine if FSS algorithms can segment all 200 bird species from the remaining images. The results are shown in Table 3. As whole object instance segmentation is the original problem domain for most of the compared methods, they show much better results than Table 1 and Table 2. Even with large variations from image to image, SST still achieves a competitive segmentation performance across bird images. Furthermore, fine-tuning SST with training images using OC-CCL again shows significant improvement in segmentation quality. We closely analyze the object instances that originally fail to be correctly segmented by SST and categorize them into 3 failure cases: *Oversegmentation*, where SST correctly segments out the object along with some extra neighboring backgrounds; *False Positive*, where SST falsely segments out an irrelevant object; and *False Negative*, where SST completely fails to segment anything from the picture. As shown in the bottom row of Fig. 8, without using any ground truth masks for most training images, fine-tuning with OC-CCL helps substantially mitigate these issues.

Table 4: **Segmentation results.** Model performance comparison on CelebAMask-HQ and MRBrainS.

Dataset	SegGPT	DEVA+SST	Cutie+SST	SAM 2+SST
CelebAMask-HQ	58.9 \pm 1.7	62.3 \pm 4.3	52.6 \pm 5.6	73.2 \pm 2.2
MRBrainS	42.7 \pm 2.7	46.6 \pm 4.3	51.4 \pm 2.9	52.6 \pm 4.1

Trait-based retrieval. As mentioned in Section 3.3, given a target trait, our method can use the reconstruction IoU to find images with similar traits. Fig. 9b shows that SST + OC-CCL faithfully retrieves images with similar corresponding traits, which can be useful for studying similar subspecies. For more experimental results and discussions, please see the Appendix.

Further exploration. SST leverages the inherent structural dependency between static specimen images to enable tracking-based segmentation. Similar dependencies exist in other domains, such as facial and medical images. Table 4 suggests superior segmentation performance of SST over the other FSS methods, indicating that SST successfully captures the underlying relationships that previous FSS methods fail to exploit. Please refer to the Appendix for more detailed results.

5 Conclusion and Discussion

We introduce Static Segmentation by Tracking (SST), a label-efficient approach for fine-grained specimen image segmentation. By applying a tracking algorithm like SAM 2 to non-sequential specimen images, SST achieves remarkable trait segmentation using only a single labeled image. Further analysis reveals that SST extends beyond specimen images, successfully segmenting animal instances in the wild. Additionally, it enables trait-level retrieval, identifying species with similar traits and patterns.

While our main use case is specimen images, this does not imply that our scope and applicability are “limited.” First, specimens are a major resource for biologists to understand organisms, and a vast amount of specimens have yet to be digitized and analyzed. Machine learning techniques are enabling efficient data processing, saving excessive manual efforts required for dealing with specimen images. Second, in addition to specimens, object-centric images with canonical poses are a common image source in various scientific fields such as MRI and CT scans. The method can be further extended for broader application. Third, while specimen images may seem easier to handle at first glance due to their object-centric nature and plain backgrounds (compared to natural images like those in MS-COCO [38]), our experiments demonstrate that segmenting fine-grained traits from them is non-trivial, particularly in a few-shot setting. In summary, our paper contributes not only to the computer vision community (*e.g.*, by promoting a rarely studied but challenging task, providing data for benchmarking, and offering a novel approach) but also to other scientific communities (*e.g.*, by facilitating the measurement of traits).

References

- [1] La variété des heliconius. <https://www.cliniquevetodax.com/Heliconius/index.html>.
- [2] Morphbank: Biological imaging. <https://www.morphbank.net/>.
- [3] Multimedia of fish specimen and associated metadata. fish-air. <https://fishair.org>.
- [4] Fmnh field museum of natural history (zoology) fish collection. *Field Museum*. https://fmipr.fieldmuseum.org/ipt/resource?r=fmnh_fishes.
- [5] Great lakes invasives network project. <https://greatlakesinvasives.org/portal/index.php>.
- [6] University of wisconsin-madison zoological museum - fish. <http://zoology.wisc.edu/wuzm/>.
- [7] Ummz university of michigan museum of zoology, division of fishes. https://ipt.lsa.umich.edu/resource?r=ummz_fish.
- [8] idigbio. <http://www.idigbio.org/portal>, 2020.
- [9] Inhs collections data. <http://biocoll.inhs.illinois.edu/portal/index.php>, 2022.
- [10] Jfbm bell atlas. <http://bellatlas.umn.edu/index.php>, 2022.
- [11] Daniel I Bolnick, Priyanga Amarasekare, Márcio S Araújo, Reinhard Bürger, Jonathan M Levine, Mark Novak, Volker HW Rudolf, Sebastian J Schreiber, Mark C Urban, and David A Vasseur. Why intraspecific trait variation matters in community ecology. *Trends in Ecology & Evolution*, 2011.
- [12] Ramakant Chandrakar, Rohit Raja, and Rohit Miri. Animal detection based on deep convolutional neural networks with genetic segmentation. *Multimedia Tools and Applications*, 2022.
- [13] Yun-Chun Chen, Yen-Yu Lin, Ming-Hsuan Yang, and Jia-Bin Huang. Show, match and segment: Joint weakly supervised learning of semantic matching and object co-segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [14] Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2022.
- [15] Ho Kei Cheng, Seoung Wug Oh, Brian Price, Alexander Schwing, and Joon-Young Lee. Tracking anything with decoupled video segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1316–1326, 2023.
- [16] Ho Kei Cheng, Seoung Wug Oh, Brian Price, Joon-Young Lee, and Alexander Schwing. Putting the object back into video object segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3151–3161, 2024.
- [17] Minsu Cho, Suha Kwak, Cordelia Schmid, and Jean Ponce. Unsupervised object discovery and localization in the wild: Part-based matching with bottom-up region proposals. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015.
- [18] Johnson N Daly M. Ohio state university fish division (osum). *Museum of Biological Diversity, The Ohio State University. Occurrence dataset*, <https://doi.org/10.15468/subsl8>, 2018.
- [19] Charles Darwin. *On the Origin of Species by Means of Natural Selection*. John Murray, London, 1859.
- [20] Richard C Edmunds, Baofeng Su, James P Balhoff, B Frank Eames, Wasila M Dahdul, Hilmar Lapp, John G Lundberg, Todd J Vision, Rex A Dunham, Paula M Mabee, et al. Phenoscape: identifying candidate genes for evolutionary phenotypes. *Molecular Biology and Evolution*, 2015.
- [21] Di Feng, Christian Haase-Schütz, Lars Rosenbaum, Heinz Hertlein, Claudius Glaeser, Fabian Timm, Werner Wiesbeck, and Klaus Dietmayer. Deep multi-modal object detection and semantic segmentation for autonomous driving: Datasets, methods, and challenges. *IEEE Transactions on Intelligent Transportation Systems*, 2020.
- [22] Isadora E. Fluck, Benjamin Baiser, Riley Wolcheski, Isha Chinniah, and Sydne Record. 2018 neon ethanol-preserved ground beetles. <https://huggingface.co/datasets/imageomics/2018-NEON-beetles>, 2024.

- [23] Abdul Mueed Hafiz and Ghulam Mohiuddin Bhat. A survey on instance segmentation: State of the art. *International Journal of Multimedia Information Retrieval*, 2020.
- [24] Sunghwan Hong, Seokju Cho, Jisu Nam, and Seungryong Kim. Cost aggregation is all you need for few-shot segmentation. *arXiv preprint arXiv:2112.11685*, 2021.
- [25] Sunghwan Hong, Seokju Cho, Jisu Nam, and Seungryong Kim. Cost aggregation is all you need for few-shot segmentation. *arXiv preprint arXiv:2112.11685*, 2021.
- [26] Grant Van Horn and macaodha. iNat challenge 2021 - FGVC8. <https://kaggle.com/competitions/inaturalist-2021>, 2021.
- [27] Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022.
- [28] Wei-Chih Hung, Varun Jampani, Sifei Liu, Pavlo Molchanov, Ming-Hsuan Yang, and Jan Kautz. Scops: Self-supervised co-part segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [29] Allan Jabri, Andrew Owens, and Alexei Efros. Space-time correspondence as a contrastive random walk. *Advances in neural information processing systems*, 33:19545–19560, 2020.
- [30] Chris Jiggins and Ian Warren. Cambridge butterfly wing collection - Chris Jiggins 2001/2 broods batch 1, 2019.
- [31] Chris Jiggins and Ian Warren. Cambridge butterfly wing collection - Chris Jiggins 2001/2 broods batch 2, 2019.
- [32] Chris Jiggins, Gabriela Montejó-Kovacevich, Ian Warren, and Eva Wiltshire. Cambridge butterfly wing collection batch 3, 2019.
- [33] Glenn Jocher, Qiu Jing, and Ayush Chaurasia. Ultralytics yolo. <https://github.com/ultralytics/ultralytics>.
- [34] Armand Joulin, Francis Bach, and Jean Ponce. Discriminative clustering for image co-segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2010.
- [35] Christopher Lawrence, Elizabeth G. Campolongo, and Neil Rosser. Heliconius collection (cambridge butterfly), 2024.
- [36] Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. Maskgan: Towards diverse and interactive facial image manipulation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [37] Xiangtai Li, Henghui Ding, Haobo Yuan, Wenwei Zhang, Jiangmiao Pang, Guangliang Cheng, Kai Chen, Ziwei Liu, and Chen Change Loy. Transformer-based visual segmentation: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- [38] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Proceedings of the European Conference of Computer Vision*, 2014.
- [39] 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.
- [40] Weide Liu, Chi Zhang, Guosheng Lin, and Fayao Liu. Crnet: Cross-reference networks for few-shot segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020.
- [41] Xiangbin Liu, Liping Song, Shuai Liu, and Yudong Zhang. A review of deep-learning-based medical image segmentation methods. *Sustainability*, 2021.
- [42] Wenhan Luo, Junliang Xing, Anton Milan, Xiaoqin Zhang, Wei Liu, and Tae-Kyun Kim. Multiple object tracking: A literature review. *Artificial Intelligence*, 2021.
- [43] Paula Mabee, James P Balhoff, Wasila M Dahdul, Hilmar Lapp, Peter E Midford, Todd J Vision, and Monte Westerfield. 500,000 fish phenotypes: The new informatics landscape for evolutionary and developmental biology of the vertebrate skeleton. *Journal of Applied Ichthyology*, 2012.

- [44] Paula M Mabee, Wasila M Dahdul, James P Balhoff, Hilmar Lapp, Prashanti Manda, Josef Uyeda, Todd Vision, and Monte Westerfield. Phenoscope: semantic analysis of organismal traits and genes yields insights in evolutionary biology. In *Application of Semantic Technology in Biodiversity Science*. IOS Press, 2018.
- [45] Anniina Mattila, Chris Jiggins, and Ian Warren. University of Helsinki butterfly wing collection - Anniina Mattila field caught specimens, 2019.
- [46] Anniina Mattila, Chris Jiggins, and Ian Warren. University of Helsinki butterfly collection - Anniina Mattila bred specimens, 2019.
- [47] Kazi Sajeed Mehrab, M Maruf, Arka Daw, Abhilash Neog, Harish Babu Manogaran, Mridul Khurana, Zhenyang Feng, Bahadir Altintas, Yasin Bakis, Elizabeth G Campolongo, et al. Fish-vista: A multi-purpose dataset for understanding & identification of traits from images. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 24275–24285, 2025.
- [48] Joana I. Meier, Patricio Salazar, Gabriela Montejó-Kovacevich, Ian Warren, and Chris Jiggins. Cambridge butterfly wing collection - Patricio Salazar PhD wild specimens batch 3, 2020.
- [49] Adriënne M Mendrik, Koen L Vincken, Hugo J Kuijff, Marcel Breeuwer, Willem H Bouvy, Jeroen De Bresser, Amir Alansary, Marleen De Bruijne, Aaron Carass, Ayman El-Baz, et al. Mrbrains challenge: online evaluation framework for brain image segmentation in 3t mri scans. *Computational intelligence and neuroscience*, 2015(1):813696, 2015.
- [50] Shervin Minaee, Yuri Boykov, Fatih Porikli, Antonio Plaza, Nasser Kehtarnavaz, and Demetri Terzopoulos. Image segmentation using deep learning: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [51] Gabriela Montejó-Kovacevich, Letitia Cookson, Eva van der Heijden, Ian Warren, David P. Edwards, and Chris Jiggins. Cambridge butterfly collection - loreto, peru 2018, 2019.
- [52] Gabriela Montejó-Kovacevich, Chris Jiggins, and Ian Warren. Cambridge butterfly wing collection batch 2, 2019.
- [53] Gabriela Montejó-Kovacevich, Chris Jiggins, and Ian Warren. Cambridge butterfly wing collection batch 4, 2019.
- [54] Gabriela Montejó-Kovacevich, Chris Jiggins, and Ian Warren. Cambridge butterfly wing collection batch 1- version 2, 2019.
- [55] Gabriela Montejó-Kovacevich, Chris Jiggins, Ian Warren, Camilo Salazar, Marianne Elias, Imogen Gavins, Eva Wiltshire, Stephen Montgomery, and Owen McMillan. Cambridge and collaborators butterfly wing collection batch 10, 2019.
- [56] Gabriela Montejó-Kovacevich, Chris Jiggins, Ian Warren, and Eva Wiltshire. Cambridge butterfly wing collection batch 5, 2019.
- [57] Gabriela Montejó-Kovacevich, Chris Jiggins, Ian Warren, and Eva Wiltshire. Cambridge butterfly wing collection batch 6, 2019.
- [58] Gabriela Montejó-Kovacevich, Chris Jiggins, Ian Warren, and Eva Wiltshire. Cambridge butterfly wing collection batch 7, 2019.
- [59] Gabriela Montejó-Kovacevich, Chris Jiggins, Ian Warren, and Eva Wiltshire. Cambridge butterfly wing collection batch 8, 2019.
- [60] Gabriela Montejó-Kovacevich, Chris Jiggins, Ian Warren, Eva Wiltshire, and Imogen Gavins. Cambridge butterfly wing collection batch 9, 2019.
- [61] Gabriela Montejó-Kovacevich, Letitia Cookson, Eva van der Heijden, Ian Warren, David P. Edwards, and Chris Jiggins. Cambridge butterfly collection - Loreto, Peru 2018 batch2, 2020.
- [62] Gabriela Montejó-Kovacevich, Letitia Cookson, Eva van der Heijden, Ian Warren, David P. Edwards, and Chris Jiggins. Cambridge butterfly collection - Loreto, Peru 2018 batch3, 2020.
- [63] Gabriela Montejó-Kovacevich, Eva van der Heijden, and Chris Jiggins. Cambridge butterfly collection - GMK Broods Ikiam 2018, 2020.
- [64] Gabriela Montejó-Kovacevich, Eva van der Heijden, Nicola Nadeau, and Chris Jiggins. Cambridge butterfly wing collection batch 10, 2020.

- [65] Gabriela Montejó-Kovacevich, Quentin Paynter, and Amin Ghane. *Heliconius erato cyrba*, Cook Islands (New Zealand) 2016, 2019, 2021, 2021.
- [66] Tai-Yu Pan, Qing Liu, Wei-Lun Chao, and Brian Price. Towards open-world segmentation of parts. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2023.
- [67] Bohao Peng, Zhuotao Tian, Xiaoyang Wu, Chengyao Wang, Shu Liu, Jingyong Su, and Jiaya Jia. Hierarchical dense correlation distillation for few-shot segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23641–23651, 2023.
- [68] Erika Pinheiro de Castro, Christopher Jiggins, Karina Lucas da Silva-Brandão, Andre Victor Lucci Freitas, Marcio Zikan Cardoso, Eva Van Der Heijden, Joana Meier, and Ian Warren. Brazilian Butterflies Collected December 2020 to January 2021, 2022.
- [69] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024.
- [70] Carsten Rother, Tom Minka, Andrew Blake, and Vladimir Kolmogorov. Cosegmentation of image pairs by histogram matching-incorporating a global constraint into mrfs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2006.
- [71] Camilo Salazar, Gabriela Montejó-Kovacevich, Chris Jiggins, Ian Warren, and Imogen Gavins. Camilo Salazar and Cambridge butterfly wing collection batch 1, 2019.
- [72] Patricio Salazar, Gabriela Montejó-Kovacevich, Ian Warren, and Chris Jiggins. Cambridge butterfly wing collection - Patricio Salazar PhD wild and bred specimens batch 1, 2018.
- [73] Patricio Salazar, Gabriela Montejó-Kovacevich, Ian Warren, and Chris Jiggins. Cambridge butterfly wing collection - Patricio Salazar PhD wild and bred specimens batch 2, 2019.
- [74] Patricio A. Salazar, Nicola Nadeau, Gabriela Montejó-Kovacevich, and Chris Jiggins. Sheffield butterfly wing collection - Patricio Salazar, Nicola Nadeau, Ikiam broods batch 1 and 2, 2020.
- [75] Zhuotao Tian, Hengshuang Zhao, Michelle Shu, Zhicheng Yang, Ruiyu Li, and Jiaya Jia. Prior guided feature enrichment network for few-shot segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [76] Cyrille Violle, Brian J. Enquist, Brian J. McGill, Lin Jiang, Céline H. Albert, Catherine Hulshof, Vincent Jung, and Julie Messier. The return of the variance: intraspecific variability in community ecology. *Trends in Ecology & Evolution*, 2012.
- [77] Huy V Vo, Patrick Pérez, and Jean Ponce. Toward unsupervised, multi-object discovery in large-scale image collections. In *Proceedings of the European Conference of Computer Vision*. Springer, 2020.
- [78] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. 2011.
- [79] Kaixin Wang, Jun Hao Liew, Yingtian Zou, Daquan Zhou, and Jiashi Feng. Panet: Few-shot image semantic segmentation with prototype alignment. In *Proceedings of the IEEE International Conference on Computer Vision*, 2019.
- [80] Xinlong Wang, Xiaosong Zhang, Yue Cao, Wen Wang, Chunhua Shen, and Tiejun Huang. Seggpt: Towards segmenting everything in context. In *Proceedings of the IEEE International Conference on Computer Vision*, 2023.
- [81] Ian Warren and Chris Jiggins. Miscellaneous *Heliconius* wing photographs (2001-2019) Part 1, 2019.
- [82] Ian Warren and Chris Jiggins. Miscellaneous *Heliconius* wing photographs (2001-2019) Part 2, 2019.
- [83] Ian Warren and Chris Jiggins. Miscellaneous *Heliconius* wing photographs (2001-2019) Part 3, 2019.
- [84] Rui Yao, Guosheng Lin, Shixiong Xia, Jiaqi Zhao, and Yong Zhou. Video object segmentation and tracking: A survey. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 11(4):1–47, 2020.
- [85] Alper Yilmaz, Omar Javed, and Mubarak Shah. Object tracking: A survey. *Acm computing surveys*, 2006.
- [86] Chi Zhang, Guosheng Lin, Fayao Liu, Rui Yao, and Chunhua Shen. Canet: Class-agnostic segmentation networks with iterative refinement and attentive few-shot learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.

Static Segmentation by Tracking: **A Label-Efficient Approach for Fine-Grained Specimen Image Segmentation**

Appendix

Table A1: **Dataset statistics** for Butterfly [35], Fish [47], and Beetle [22] dataset.

	Butterfly		Fish	Beetle
	Major	Minor		
# of Classes	5	146	474	12
Total Train	2,831	-	-	120
Total Test	500	313	1,573	60

A1 Dataset Statistics and Evaluation Details

We put the general statistics for Butterfly [35], Fish [47], and Beetle [22] datasets in Table A1.

A1.1 Butterfly.

The Cambridge Butterfly [35] dataset includes 151 butterfly subspecies of the *Heliconius* genus; each has 4 ~ 14 distinctive traits to tell itself apart from others. Example traits include the tiger tails on the hindwings and white bands on the forewings; some have quite complex, disconnected shapes. Across specimens of the same subspecies, the mask IDs are consistent. An algorithm needs to segment them and also label each with an ID.

We split the dataset into two parts, Major and Minor, based on the available data sample for each subspecies. The Major part has 5 different subspecies, with 2,831 semi-automatically labeled training samples and 500 hand-labeled test samples in total. The semi-automatic approach: we used SST to propagate masks, followed by human inspection. *Samples with unsatisfactory masks were then hand-labeled.* We take all 2,831 training data samples to train standard segmentation models, Mask2Former [14] and YOLOv8 [33], for each subspecies. To evaluate few-shot segmentation models, we sample a specified number of random specimens from the train set for each subspecies, and evaluate the performance on the corresponding test data.

The Minor part has 146 subspecies with 313 hand-labeled test images in total. As there is insufficient data for each subspecies to construct a training set, Minor is intended for the one-shot segmentation task. For this part, we only test on few-shot segmentation models in the same fashion as we evaluate the Major subspecies.

A1.2 Fish.

The Fish-Vista [47] dataset has 474 different species of fish, containing 1,573 samples in total. All species share a common set of 9 segmentation classes (*e.g.*, head, eye, tail, adipose fin, caudal fin, etc.). As there are a limited number of samples per species, we do not run many-shot model training or OC-CCL fine-tuning on Fish. For few-shot segmentation, we use 1 sample from each species as reference, and predict the segmentation masks of the other samples.

A1.3 Beetle.

The Beetle [22] dataset consists of beetles of 12 different species. Each beetle species shares 5 common segmentation classes: head, pronotum, elytra, antenna, and legs. The antennae and leg parts of beetles are quite challenging, with complex, sharp, and thin shapes. As the beetle species share similar visual traits, we always apply a universal model across all the species. For each species, we hand-label 15 images in total, taking 10 as the training set and 5 as the test set. For the many-shot instance segmentation models, we train one single model on all training data across species. For

few-shot segmentation methods and SST, we sample one example from the 120 training samples and test the segmentation quality on all 60 test samples across species.

A1.4 Remark.

We emphasize that while we used SAM 2 to assist in data annotation, *humans* inspected the results and re-labeled them when necessary. Additionally, the test set was entirely labeled by *humans*. Consequently, a well-trained model, such as Mask2Former [14], given sufficient labeled data, could surpass SST in performance.

A2 Implementation Details

A2.1 Many-Shot Methods

Existing many-shot methods usually rely on sufficient training data to obtain ideal segmentation performance. As shown in Fig. 3, YOLOv8 [33] and Mask2Former [14] require more than 1,000 training samples to achieve comparable performance with SST. That means the many-shot methods cannot work as expected for Butterfly Minor [35] and Fish [47], where insufficient data is provided within each (sub)species. Therefore, we only train many-shot models for Butterfly Major [35] and Beetle [22] datasets. More specifically, for Butterfly Major [35], we train independent segmentation models for each subspecies, while for Beetle [22], only one model is trained for the whole dataset. The training is conducted following standard hyper-parameter settings of the adopted methods.

A2.2 OC-CCL Fine-tuning

Similar to many-shot learning methods, OC-CCL fine-tuning cannot be applied to Butterfly Minor [35] or Fish [47] due to insufficient samples with each (sub)species. For Beetle [22] and Butterfly Major [35], we follow the default LoRA initialization and learning rate scheduling for LoRA [27], and train it for one epoch. More specifically, we replace all linear projection layers and MLP layers in the memory encoder and mask decoder of the video segmentation models with LoRA linear layers, using $r = 32$ and $\alpha = 64$. We take one fixed labeled image as x_0 , and iteratively sample x_1 from the unlabeled training set to create multiple short palindrome cycles, as introduced in Section 3.2. For Beetle [22], due to the inherently high variation in beetle orientations, we introduce a small degree of random 2D rotation augmentation on the labeled example during the fine-tuning stage to enhance the robustness of the model towards rotation OOD scenarios.

Table A2: **OC-CCL fine-tuning time** on different datasets for one epoch (minutes).

Model	Major [35]	Beetle [22]
DEVA [15]	10	5
Cutie [16]	10	4
SAM 2 [69]	6	2

As illustrated in Table A2, the fine-tuning does not require excessive time for the adopted datasets. It enhances the application value for SST on real-world scenarios.

A2.3 One-Shot Segmentation Setting

For the Butterfly Major [35] and Beetle [22] datasets, we perform 20 runs of one-shot segmentation experiments to ensure a fair and stable comparison. In each run, we sample a different labeled image from the training set as the reference frame and evaluate across all testing set examples. The average mIoU is then calculated across all runs with the standard deviation reported.

For Fish [47] and Butterfly Minor [35], where there is insufficient data to sample 20 reference images, we conduct only one single run. Specifically, we use one labeled example as the reference and perform inference on the remaining data within the (sub)species.

A2.4 Multi-Shot Setting

Similar to the other few-shot segmentation algorithms, SST can also benefit under the multi-shot setting. Given more labeled images, video segmentation models leverage them as multiple reference frames to better capture information for the subsequent prediction steps. For 5-shot evaluation, we adopt a naive design, where five labeled image mask pairs are put in the beginning of the video sequence, where a new sequence is created independently for each of the incoming images. In other words, given labeled reference pairs $(x_0, y_0), (x_1, y_1), (x_2, y_2), (x_3, y_3), (x_4, y_4)$ and unlabeled target images $\{x_5, \dots, x_N\}$, we create sequences $\{x_0, \dots, x_4, x_5\}, \dots, \{x_0, \dots, x_4, x_N\}$ and apply SST to each of them. We also aim to design a better multi-shot mechanism for SST as a future work.

A3 More Baselines

In addition to FSS and many-shot instance segmentation models, we also adopt a naive “copy and paste” method and optical flow as supplementary baselines.

Table A3: **Segmentation results** across different datasets using copy baseline and optical flow.

	Major [35]	Minor [35]	Fish [47]	Beetle [22]
Copy	20.9 \pm 10.9	36.1	10.6	13.2 \pm 6.7
Optical Flow	24.9 \pm 11.9	36.2	5.7	16.5 \pm 4.8
DEVA [15] + SST	73.1 \pm 4.0	68.6	50.8	39.0 \pm 6.9
Cutie [16] + SST	67.4 \pm 5.0	69.7	51.9	45.8 \pm 4.4
SAM 2 [69] + SST	81.0 \pm 1.0	70.6	70.4	61.9 \pm 3.7

A3.1 Copy Baseline.

Given the highly aligned nature, an intuitive way to address the fine-grained specimen image segmentation problem is simply copying the segmentation masks from the reference image. We evaluate the copy baseline in Table A3. While the naive solution achieves certain mIoU performance by chance, we demonstrate that on all three video segmentation models, SST yields much stronger results. It also indicates that although the problem seems straightforward, it takes non-trivial efforts to achieve practical results.

A3.2 Optical Flow.

Another intuitive solution to this problem is optical flow, as it’s a fundamental way to calculate the pixel mapping from one image to the other. This technique is widely used in the medical image registration and video tracking field, which is similar to our task. To implement this, we first perform dense optical flow to calculate the per-pixel mapping from the labeled to unlabeled images. We then send the segmentation masks through the same mapping to “map” the masks to the corresponding region on the target image. As shown in Table A3, there is limited improvement using optical flow compared with the copy baseline. It further demonstrates that specimen trait segmentation cannot be easily solved through existing techniques.

A4 Discussion.

A4.1 Fairness of Different Baselines.

We note that different compared methods may vary in their parameters, pre-trained data, training and inference time, etc. Aligning them completely is challenging, especially since SST uses video models, which fundamentally differs from an image model in several aspects. That said, we use video models not for their training data, GPUs, or size, but for their video segmentation capability, which unexpectedly aligns with our problem.

A4.2 Images with Multiple Specimens.

We assume that each image contains a single specimen instance. If an image contains multiple specimens, object detectors like Grounding DINO [39] can be applied beforehand to separate them.

A4.3 Effectiveness of Video Segmentation Models.

We view this as an emergent property of models trained for video segmentation—they can track and segment taxonomically related species across non-sequential, independently captured photographs.

A5 Additional Analysis on SST

A5.1 Computational Cost

The inference stage of SST involves mask propagation across images, which might raise concerns about the computational cost. In Table 1, we report the time each method requires to process a single instance on one NVIDIA A100 GPU with 40 GB of VRAM. SST demonstrates greater efficiency than most of the compared methods, confirming its practical applicability.

A5.2 Detailed Results on Face and MRI Images

Table A4: **Segmentation results** of different models across facial features on CelebAMask-HQ [36].

	Overall	Eyes	Nose	Mouth
SegGPT [80]	58.9 \pm 1.7	46.7 \pm 2.6	68.7 \pm 3.3	61.3 \pm 5.7
DEVA [15] + SST	62.3 \pm 4.3	53.2 \pm 5.6	64.1 \pm 8.3	69.7 \pm 5.5
Cutie [16] + SST	52.6 \pm 5.6	54.6 \pm 6.6	46.9 \pm 14.2	56.4 \pm 5.0
SAM 2 [69] + SST	73.2 \pm 2.2	70.3 \pm 1.4	64.2 \pm 5.3	81.8 \pm 2.2

Table A5: **Segmentation results** of different models across brain tissue types on MRBrainS [49].

	Overall	Cerebrospinal Fluid	Gray Matter	White Matter
SegGPT [80]	42.7 \pm 2.7	39.1 \pm 1.7	45.6 \pm 1.1	43.4 \pm 1.2
DEVA [15] + SST	46.6 \pm 4.3	42.0 \pm 3.7	52.3 \pm 2.1	45.6 \pm 2.0
Cutie [16] + SST	51.4 \pm 2.9	49.8 \pm 1.9	55.5 \pm 2.5	48.8 \pm 2.4
SAM 2 [69] + SST	52.6 \pm 4.1	48.1 \pm 3.1	51.6 \pm 3.5	58.1 \pm 1.9

We report the detailed results on CelebAMask-HQ [36] and MRBrainS [49] for different segmentation targets in Table A5 and Table A4, respectively. SST yields better performance than SegGPT [80] on most of the tasks.

A5.3 Analysis on Inference Variants

In the main paper, we evaluate our method using a single test image at a time to compare it against other few-shot segmentation models, ensuring a fair comparison. However, as mentioned in Section 3.3, it is also possible to concatenate all test images and process them together using SST. We observe that using random orders does not significantly degrade performance, though a more carefully designed algorithm could make SST more stable.

To demonstrate this, we conduct a toy experiment on the *lativitta* subspecies from the Butterfly [35] dataset. We first randomly sample 10 butterflies from the test set and generate 10,000 unique orderings for the 10 images. We then put all ten images in a sequence based on each ordering and evaluate SST after propagating through the entire sequence. We average the mIoU performance across each ordering and plot it against a histogram as shown in the blue part of Fig. A1. There appear to be two peaks in the distribution, but the overall influence is limited (mIoU from 90% to 92%). We then experiment with a slightly improved ordering strategy by interleaving each test image with the reference image, so that the reference information can be retained even if some frames lose the track.

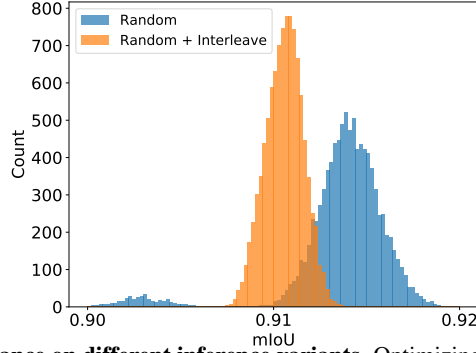


Figure A1: **SST performance on different inference variants.** Optimizing the frame ordering improves the stability when using long video sequences for inference.

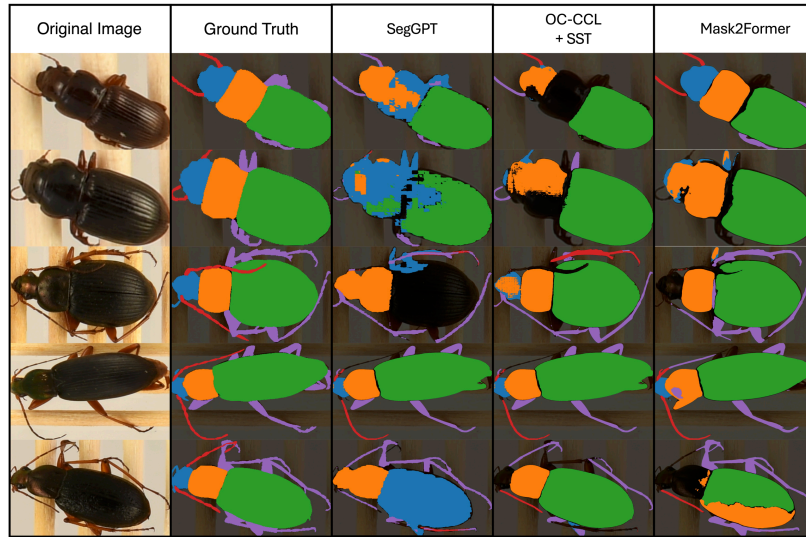


Figure A2: **Qualitative results on Beetle [22].** After fine-tuning, SST demonstrates segmentation results better than SegGPT [80] and comparable with Mask2Former [14]. SST is applied to SAM 2 [69] in this experiment.

As shown in the orange part of Fig. A1, although random ordering only has a limited influence on the performance, interleaving the test images further reduces the standard deviation in mIoU. Thus, we conclude that finding the optimal sequence can indeed help with the stability of video inference, and we plan to explore the ordering design as a future work.

A6 Qualitative Results on Fine-Tuned Models

To compare the performance of different methods, we show more qualitative results on a more diverse set of butterflies, fish, and beetle species, see Fig. A2, Fig. A3, and Fig. A4. For each column, we keep the same setting as in the main paper. For both SegGPT [80] and SST + OC-CCL, we select one random image from the training set as a reference and evaluate the segmentation quality on the target image. The quality is demonstrated in the third and fourth columns of the figures. We also show the segmentation quality of Mask2Former [14] in the last column of Fig. A2 and Fig. A4, which is trained on the entire available training dataset. We omit the Mask2Former column for the Butterfly [35] dataset as there aren't enough data samples to train a full standard segmentation model for most of these subspecies.

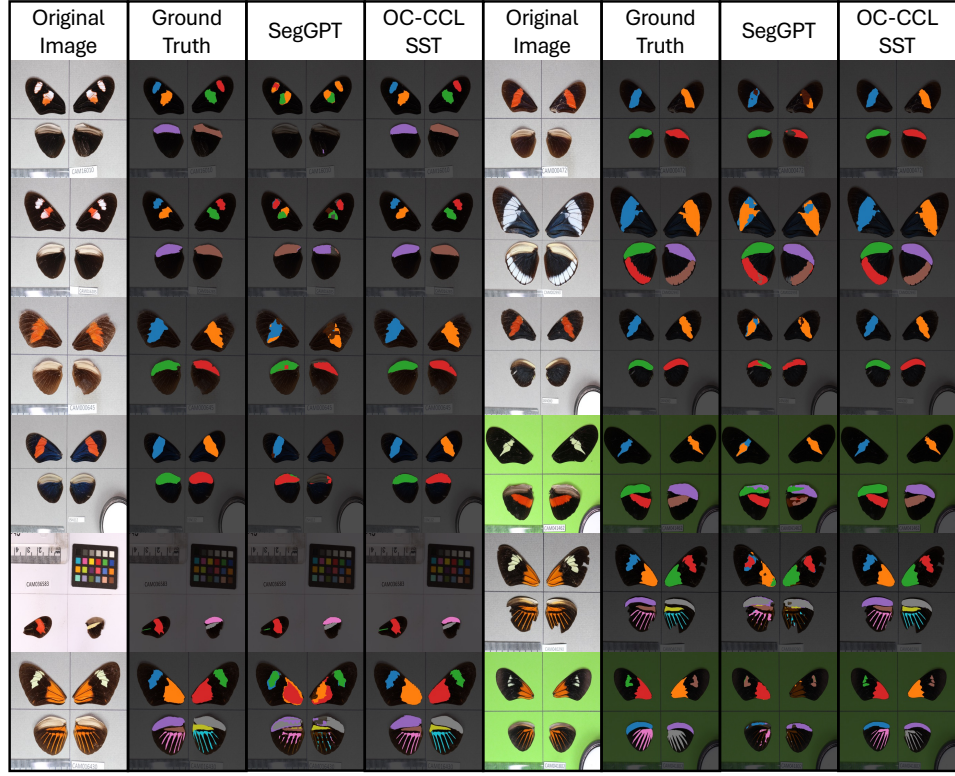


Figure A3: **Qualitative results on Butterfly [35]**. SST and OC-CCL are applied to SAM 2 [69].

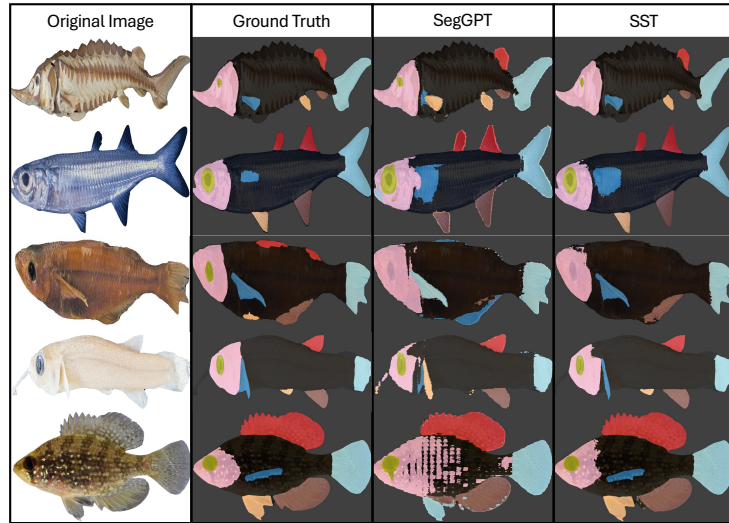


Figure A4: **Qualitative results on Fish [47]**. SST and OC-CCL are applied to SAM 2 [69].

A7 Additional Trait-Based Retrieval Results

We originally demonstrated SST’s ability to do trait-based retrieval in Sections 3.4 and 4.5 in the main paper. Here, we show more trait-based retrieval results using SST with a diverse range of subspecies. Given a target trait on any subspecies, as outlined in red in the left-most column of Fig. A5, SST can reliably retrieve subspecies that share similar traits, as outlined in cyan in Fig. A5.

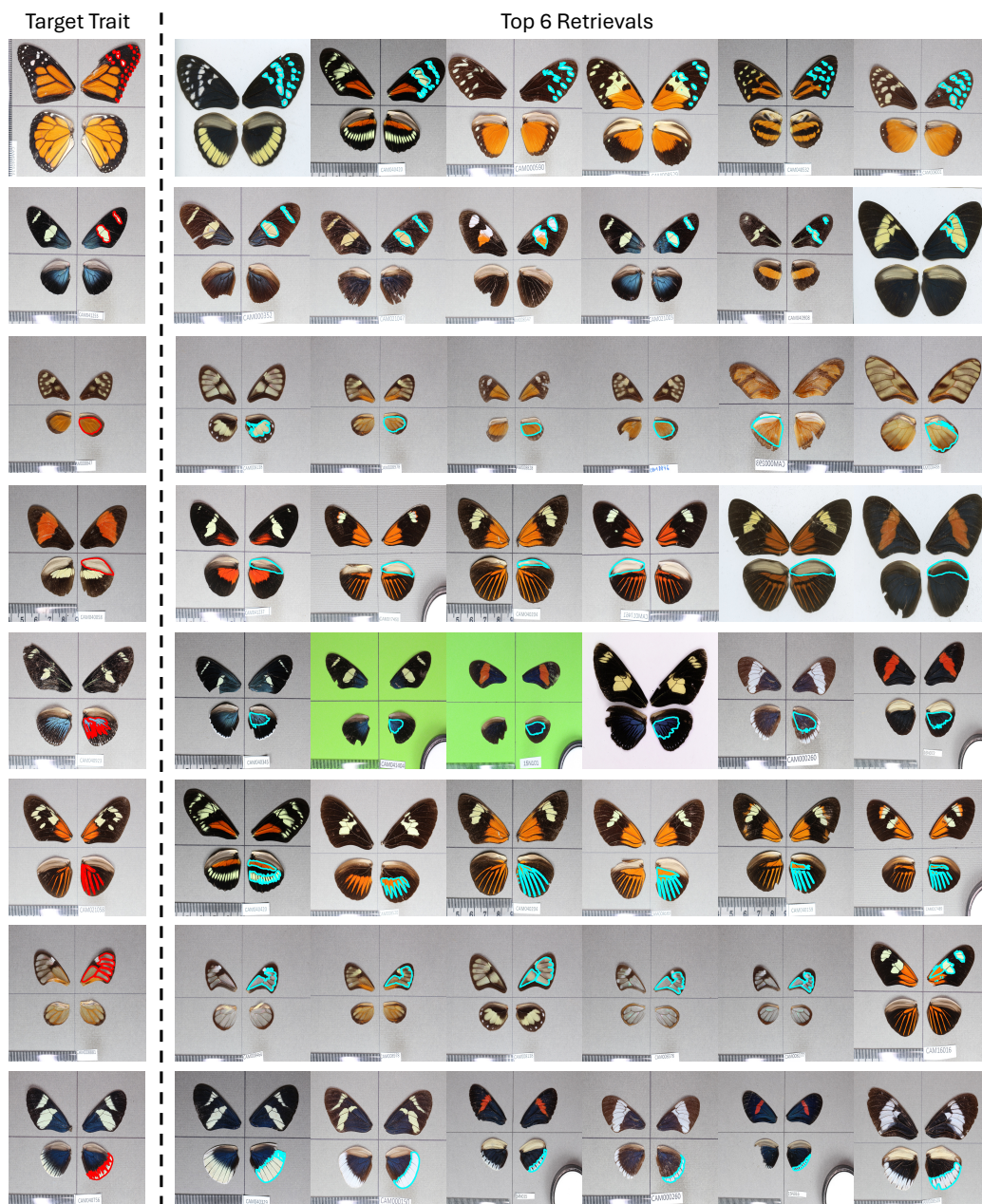


Figure A5: Qualitative results for trait retrieval on Butterfly [35]. Target trait is outlined in red, the retrieved traits are outlined in cyan.