

# ZERO-SHOT TREE DETECTION AND SEGMENTATION FROM AERIAL FOREST IMAGERY

**Michelle Chen<sup>1</sup> David Russell<sup>2</sup> Amritha Pallavoor<sup>2</sup> Derek Young<sup>2</sup> Jane Wu<sup>1</sup>**

<sup>1</sup> UC Berkeley <sup>2</sup> UC Davis

<sup>1</sup> {michelle.chenn, janehwu}@berkeley.edu

<sup>2</sup> {djrussell, aspallavoor, djyoung}@ucdavis.edu

## ABSTRACT

Large-scale delineation of individual trees from remote sensing imagery is crucial to the advancement of ecological research, particularly as climate change and other environmental factors rapidly transform forest landscapes across the world. Current RGB tree segmentation methods rely on training specialized machine learning models with labeled tree datasets. While these learning-based approaches can outperform manual data collection when accurate, the existing models still depend on training data that's hard to scale. In this paper, we investigate the efficacy of using a state-of-the-art image segmentation model, Segment Anything Model 2 (SAM2), in a *zero-shot* manner for individual tree detection and segmentation. We evaluate a pretrained SAM2 model on two tasks in this domain: (1) zero-shot segmentation and (2) zero-shot transfer by using predictions from an existing tree detection model as prompts. Our results suggest that SAM2 not only has impressive generalization capabilities, but also can form a natural synergy with specialized methods trained on in-domain labeled data. We find that applying large pretrained models to problems in remote sensing is a promising avenue for future progress. We make our code available at: <https://github.com/open-forest-observatory/tree-detection-framework>.

## 1 INTRODUCTION

Remote sensing of global forest landscapes plays a critical role in a broad range of research areas, from ecological monitoring (Lechner et al., 2020; Fassnacht et al., 2024) to wildfire prevention (Iban & Sekertekin, 2022) to climate change resiliency (Falk et al., 2022; Forzieri et al., 2022). In particular, the ability to identify and measure individual trees provides valuable insights into forest structure, biodiversity, carbon sequestration, and overall ecosystem health. Traditional methods of forest inventory, such as ground-based surveys, are time-consuming, labor-intensive, and often limited in scale. In contrast, remotely sensed overhead imagery (e.g. via unmanned aerial vehicles) combined with machine learning techniques offers a promising solution for large-scale tree analysis.

Existing machine learning methods for tree detection are CNN-based models that were trained or fine-tuned on manually labeled tree crown datasets. DeepForest (Weinstein et al., 2020) uses the RetinaNet architecture (Lin, 2017) and retrains from the lidar-based method in Silva et al. (2016) using RGB imagery from sites in the National Ecological Observatory Network (NEON) dataset (Weinstein et al., 2019). In a similar vein, Detectree2 (Ball et al., 2023) introduced a Mask R-CNN model (He et al., 2017) based on Detectron2 and trained on images collected in Malaysia and French Guiana. Because these approaches rely on models that were only trained on relatively small, labeled tree datasets, they can naturally struggle with generalization to unseen geographical areas. Following the rise of transformer-based models (Vaswani, 2017) trained on Internet-scale data, there is growing evidence that large pretrained models can be effectively applied in a zero-shot manner to various scientific domains (see e.g. Cambrin et al. (2024); Gutiérrez et al. (2024)).

Building on recent advancements in computer vision, we investigate the efficacy of using a pretrained Segment Anything Model 2 (SAM2) (Ravi et al., 2024) in a zero-shot manner for automatic tree detection and segmentation from aerial imagery. Utilizing transformer-based architecture with proven zero-shot generalization capabilities provides significant benefits over previous methods, and

our proposed approach could enhance generalization across diverse ecosystems around the world. In order to evaluate the generalization capabilities of SAM2, we investigate (1) zero-shot prediction using a pretrained SAM2 model and (2) whether zero-shot transfer from tree detection to segmentation is viable by prompting SAM2 using bounding boxes predicted by DeepForest. We benchmark pretrained SAM2 against DeepForest and Detectree2 using two standard datasets from NEON and Detectree2. Furthermore, qualitative analysis of Emerald Point imagery from the Open Forest Observatory<sup>1</sup>, which previously served as the basis for geometric tree detection optimization (Young et al., 2022), demonstrates that SAM2 significantly outperforms existing models. Our experiments highlight the potential for large pretrained models in machine learning to drive progress in remote sensing, particularly for visual understanding at scale.

## 2 ZERO-SHOT PREDICTION AND TRANSFER

**Segment Anything Model 2.** SAM2 (Ravi et al., 2024) is a foundation model designed for promptable visual segmentation across images and videos. Given an input RGB image and a set of point, bounding box, or mask prompts, the model outputs segmentation masks conditioned on the prompts. The neural network architectures for both SAM2 and its predecessor, SAM (Kirillov et al., 2023), consist of an encoder-decoder structure using transformers. The image encoder is a vision transformer (Ryali et al., 2023) pretrained using Masked Autoencoders (MAE) (He et al., 2022), and the mask decoder includes a stack of transformer blocks to update the image and prompt embeddings in a bidirectional manner. In addition to performance enhancements on images, SAM2 extends SAM to the video domain to enable both segmentation and tracking across time. At present, SAM2 represents the current state-of-the-art for general image and video segmentation.

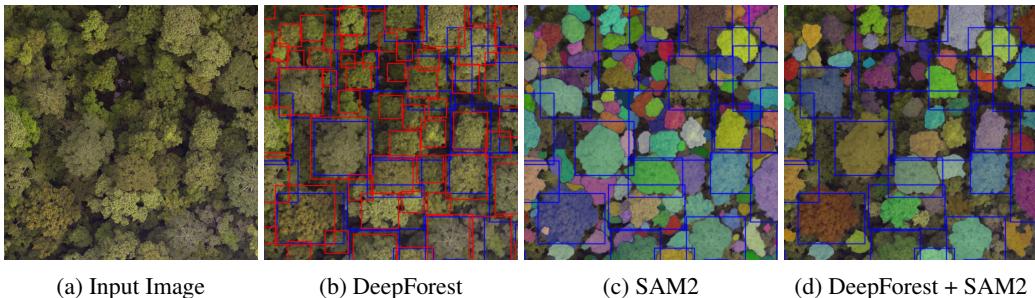
**Zero-shot Tree Segmentation.** We first aim to understand SAM2’s zero-shot capabilities as a standalone method for automatic tree segmentation and detection. Throughout our study, we use a pretrained SAM2 model with the Hiera-L image encoder and refer to this zero-shot prediction approach as “SAM2.” To segment trees across an image, we utilize SAM2’s automatic mask generator, which uniformly samples single-point prompts in a grid structure that each yields a corresponding mask. Due to the large-scale format of aerial imagery, we find that running the model on cropped sections and applying the standard post-processing to remove disconnected regions significantly improves segmentation quality. After mask generation, bounding box predictions can be calculated for each mask to obtain tree detection results, which we use for quantitative evaluation on labeled tree crown data. Example SAM2 zero-shot predictions are shown in Figure 1c.

**Zero-shot Transfer from Tree Detection.** SAM’s prompt-based segmentation approach is designed for zero-shot transfer to downstream segmentation tasks. This means that through prompt engineering, knowledge from *domain-specific models* can be directly transferred to SAM/SAM2. In the context of vision-based tree delineation, the vast majority of existing models (Fromm et al., 2019; Weinstein et al., 2020; Ball et al., 2023; Gan et al., 2023) are tree crown detection models that output either bounding boxes or polygons, which are much less labor-intensive to label than segmentation masks. Aiming to build on this domain-specific expertise, we also investigate the efficacy of providing bounding box predictions from specialized tree detectors as prompts to SAM2. Example SAM2 predictions using DeepForest bounding boxes are shown in Figure 1d.

## 3 EXPERIMENTS

**Datasets.** We evaluate SAM2’s zero-shot capabilities on images drawn from three datasets spanning two continents: the Emerald Point dataset (Young et al., 2022), the National Ecological Observatory Network (NEON) TreeEvaluation dataset (Weinstein et al., 2019; 2021), and the released Detectree2 dataset (Ball et al., 2023). First, in order to evaluate the generalization of both SAM2 and existing models, we use a dataset of aerial photographs captured by a quadcopter (referred to as Emerald Point) collected at a study site in Emerald Bay State Park located in California (Young et al., 2022). For quantitative benchmarking, we use the NEON TreeEvaluation dataset, which includes 30,975 manually-annotated tree bounding boxes in 22 sites across the United States. We also use the Danum, Sepilok East, and Sepilok West datasets released as part of Detectree2, which includes approximately 1,000 labeled crowns from two tropical field sites in Malaysia.

<sup>1</sup><https://openforestobservatory.org>



(a) Input Image (b) DeepForest (c) SAM2 (d) DeepForest + SAM2

Figure 1: SAM’s prompt-based segmentation framework enables (c) *zero-shot tree segmentation* and (d) *zero-shot transfer* of tree detection models such as (b) DeepForest to SAM2 via bounding box prompting. An example image from the Detectree2 dataset is used. Ground truth bounding boxes are drawn in blue, predicted bounding boxes are in red, and predicted masks are randomly colored.

**Baselines.** We compare SAM2 with two CNN-based models for image-based tree detection: DeepForest (Weinstein et al., 2020) and Detectree2 (Ball et al., 2023). DeepForest was trained on NEON data, and Detectree2 was trained on the Danum, Sepilok, and Paracou sites in their collected dataset. The specific models we use are the DeepForest 1.5.0 release and the Detectree2 2.0.1 release, both with default parameters. For all methods, including SAM2, we apply non-maximum suppression (NMS) suppression to predicted polygons with an intersection-over-union (IOU) threshold of 0.05. Applying NMS directly to polygon representations retains more detections compared to using bounding boxes, resulting in higher precision but lower recall values. To the best of our knowledge, the data we evaluate all methods on was not used during training.

**Generalization to Emerald Point.** Neither DeepForest/Detectree2 nor SAM2 was trained using data from Emerald Bay State Park, so we can evaluate how well all three models generalize to unseen images/geographical regions using the Emerald Point dataset. Figure 2 shows qualitative comparisons between predicted DeepForest bounding boxes, Detectree2 polygons, and SAM2 polygon masks generated using automatic mask generation. SAM2 is able to capture smaller and shorter trees with high accuracy, offering more detailed and reliable segmentation across varied forest densities. Though DeepForest outputs bounding boxes and Detectree2 outputs polygons, the number and quality of detected trees are significantly lower than what SAM2 predicts. The drop in performance of existing models could be due in part to a change in camera perspective, as existing approaches were trained only on top-down orthomosaics. In contrast, SAM2 demonstrates the ability to effectively perform zero-shot segmentation from oblique perspectives in addition to top-down views. While more extensive evaluations are needed, our results are an early indication that a pretrained SAM2 already shows impressive generalization capabilities to tree segmentation.

**Benchmarking on tree crown datasets.** While SAM2 is not technically an object detection model, we nonetheless provide some quantitative evaluation of its tree delineation capabilities by benchmarking on the NEON TreeEvaluation and Detectree2 datasets. Table 1 compares SAM2’s zero-shot tree detection with DeepForest and Detectree2, both of which are performing in-distribution inference for their corresponding datasets (NEON and Detectree2). For each dataset, we compute average precision and recall at an IOU threshold of 0.4 to evaluate the quality of tree crown bounding box predictions, with a minimum confidence for detections threshold of 0.1 (note that this does not exactly match how prior work metrics are reported). Because SAM2 is (1) not meaningfully prompted and (2) does not distinguish between object classes, we observe that the model tends to over segment, which results in the low precision scores. However, the results also indicate that SAM2 can achieve comparable recall to existing methods, particularly Detectree2, across both datasets. Figure 3 shows representative results from both the NEON and Detectree2 datasets.

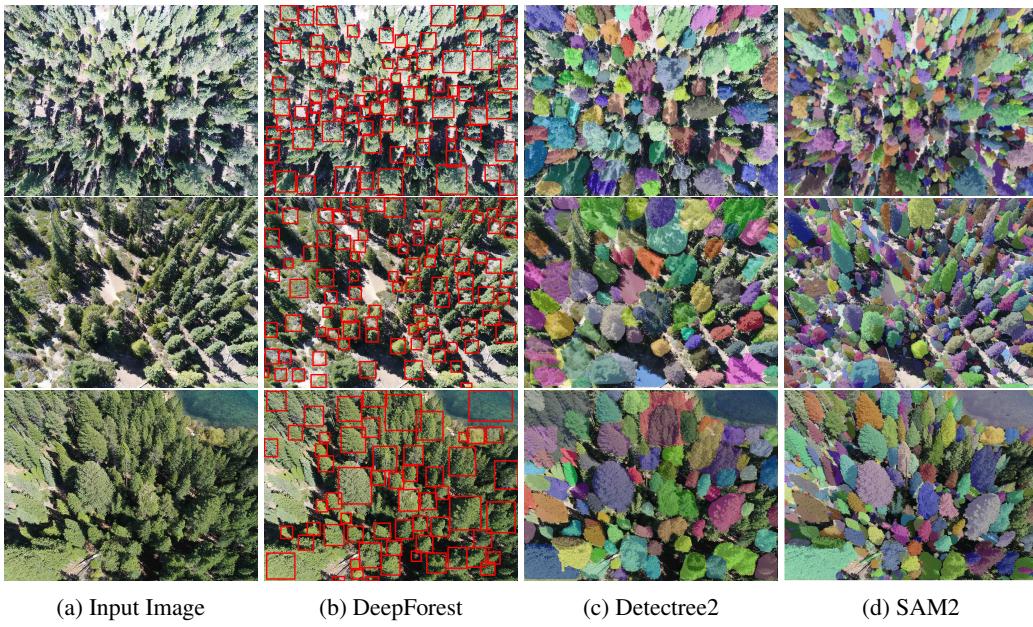


Figure 2: Segmentation results using aerial images from the Emerald Point dataset. Each row includes the input image, DeepForest tree crown bounding boxes, Detectree2 polygons, and SAM2 segmentation masks.

Method	NEON		Detectree2	
	Precision ( $\uparrow$ )	Recall ( $\uparrow$ )	Precision ( $\uparrow$ )	Recall ( $\uparrow$ )
DeepForest	<b>0.61</b>	<b>0.76</b>	0.34	0.60
Detectree2	0.25	0.50	<b>0.40</b>	<b>0.72</b>
SAM2	0.17	0.41	0.16	<u>0.65</u>

Table 1: Evaluation of tree crown bounding box detection on the NEON and Detectree2 benchmark datasets using DeepForest, Detectree2, and a pretrained SAM2 model (zero-shot). For all methods, bounding box precision and recall are determined using an IoU threshold of 0.4 and computed independently for each image before averaging across all images.

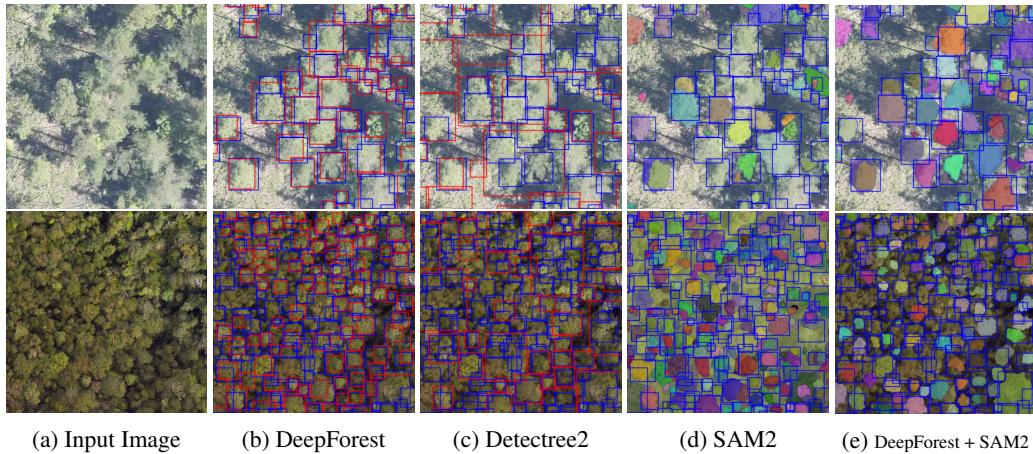


Figure 3: Tree crown detection results using RGB images from the NEON (Row 1) and Detectree2 (Row 2) datasets. Ground truth bounding boxes are drawn in blue, and predictions are drawn in red. Each row includes the input image, DeepForest, Detectree2, and SAM2, and the bounding box prompted SAM2. Note that SAM2 tends to over-segment in the absence of prompts, which contributes to the low precision values in Table 1.

## 4 DISCUSSION

In this paper, we explore the zero-shot capabilities of SAM2 for tree delineation in aerial remote sensing imagery. The existing paradigm for tree detection and segmentation is to train specialized models using labeled tree crown data; however, we conjecture that recent progress in foundation models like SAM/SAM2 can greatly enhance the generalization capabilities of vision-based techniques in remote sensing. We focus on two specific tasks in the domain of tree delineation: (1) zero-shot tree segmentation using SAM2 and (2) zero-shot transfer of tree detection to SAM2 segmentation. Our experimental results show strong indications that SAM2 has the potential to generalize to a wide range of tree species, canopy structures, and environmental conditions, perhaps even without any dataset-specific adaptation or fine-tuning.

## ACKNOWLEDGEMENTS

We would like to thank Marissa Ramirez de Chanlatte, Trevor Darrell, and Arjun Rewari for helpful discussions and feedback. This work was supported in part by the National Science Foundation Division of Biological Infrastructure (#2152671, #2152672, #2152673). J. W. was supported by the NSF Mathematical Sciences Postdoctoral Fellowship and the UC President’s Postdoctoral Fellowship.

## REFERENCES

- James GC Ball, Sebastian HM Hickman, Tobias D Jackson, Xian Jing Koay, James Hirst, William Jay, Matthew Archer, Mélaine Aubry-Kientz, Grégoire Vincent, and David A Coomes. Accurate delineation of individual tree crowns in tropical forests from aerial rgb imagery using mask r-cnn. *Remote Sensing in Ecology and Conservation*, 9(5):641–655, 2023.
- Daniele Rege Cambrin, Isaac Corley, and Paolo Garza. Depth any canopy: Leveraging depth foundation models for canopy height estimation. *arXiv preprint arXiv:2408.04523*, 2024.
- Donald A Falk, Philip J van Mantgem, Jon E Keeley, Rachel M Gregg, Christopher H Guiterman, Alan J Tepley, Derek JN Young, and Laura A Marshall. Mechanisms of forest resilience. *Forest Ecology and Management*, 512:120129, 2022.
- Fabian Ewald Fassnacht, Joanne C White, Michael A Wulder, and Erik Næsset. Remote sensing in forestry: current challenges, considerations and directions. *Forestry: An International Journal of Forest Research*, 97(1):11–37, 2024.
- Giovanni Forzieri, Vasilis Dakos, Nate G McDowell, Alkama Ramdane, and Alessandro Cescatti. Emerging signals of declining forest resilience under climate change. *Nature*, 608(7923):534–539, 2022.
- Michael Fromm, Matthias Schubert, Guillermo Castilla, Julia Linke, and Greg McDermid. Automated detection of conifer seedlings in drone imagery using convolutional neural networks. *Remote Sensing*, 11(21):2585, 2019.
- Yi Gan, Quan Wang, and Atsuhiro Iio. Tree crown detection and delineation in a temperate deciduous forest from uav rgb imagery using deep learning approaches: Effects of spatial resolution and species characteristics. *Remote Sensing*, 15(3):778, 2023.
- Juan D Gutiérrez, Roberto Rodriguez-Echeverria, Emilio Delgado, Miguel Ángel Suero Rodrigo, and Fernando Sánchez-Figueroa. No more training: Sam’s zero-shot transfer capabilities for cost-efficient medical image segmentation. *IEEE Access*, 2024.
- Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pp. 2961–2969, 2017.
- Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 16000–16009, 2022.

- Muzaffer Can Iban and Alihsan Sekertekin. Machine learning based wildfire susceptibility mapping using remotely sensed fire data and gis: A case study of adana and mersin provinces, turkey. *Ecological Informatics*, 69:101647, 2022.
- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4015–4026, 2023.
- Alex M Lechner, Giles M Foody, and Doreen S Boyd. Applications in remote sensing to forest ecology and management. *One Earth*, 2(5):405–412, 2020.
- T Lin. Focal loss for dense object detection. *arXiv preprint arXiv:1708.02002*, 2017.
- 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.
- Chaitanya Ryali, Yuan-Ting Hu, Daniel Bolya, Chen Wei, Haoqi Fan, Po-Yao Huang, Vaibhav Aggarwal, Arkabandhu Chowdhury, Omid Poursaeed, Judy Hoffman, et al. Hiera: A hierarchical vision transformer without the bells-and-whistles. In *International Conference on Machine Learning*, pp. 29441–29454. PMLR, 2023.
- Carlos A Silva, Andrew T Hudak, Lee A Vierling, E Louise Loudermilk, Joseph J O’Brien, J Kevin Hiers, Steve B Jack, Carlos Gonzalez-Benecke, Heezin Lee, Michael J Falkowski, et al. Imputation of individual longleaf pine (*pinus palustris* mill.) tree attributes from field and lidar data. *Canadian journal of remote sensing*, 42(5):554–573, 2016.
- A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.
- Ben G Weinstein, Sergio Marconi, Stephanie Bohlman, Alina Zare, and Ethan White. Individual tree-crown detection in rgb imagery using semi-supervised deep learning neural networks. *Remote Sensing*, 11(11):1309, 2019.
- Ben G Weinstein, Sergio Marconi, Mélaine Aubry-Kientz, Gregoire Vincent, Henry Senyondo, and Ethan P White. Deepforest: A python package for rgb deep learning tree crown delineation. *Methods in Ecology and Evolution*, 11(12):1743–1751, 2020.
- Ben G Weinstein, Sergio Marconi, Stephanie A Bohlman, Alina Zare, Aditya Singh, Sarah J Graves, and Ethan P White. A remote sensing derived data set of 100 million individual tree crowns for the national ecological observatory network. *Elife*, 10:e62922, 2021.
- Derek JN Young, Michael J Koontz, and JonahMaria Weeks. Optimizing aerial imagery collection and processing parameters for drone-based individual tree mapping in structurally complex conifer forests. *Methods in Ecology and Evolution*, 13(7):1447–1463, 2022.