

Evaluating Deep Models for Tree Canopy Segmentation Under Extreme Data Scarcity

David Szczecina, Hudson Sun , Niloofar Azad, Anthony Bertnyk

{david.szczecina, hudson.sun, n2azad, abertnyk}@uwaterloo.ca

Abstract

Tree canopy detection from aerial imagery is an important task for environmental monitoring, urban planning, and ecosystem analysis. The Solafune Tree Canopy Detection competition [14] provides a small and imbalanced dataset of only 150 annotated images, posing significant challenges for training deep models without severe overfitting. In this work, we evaluate five representative architectures, YOLOv11, Mask R-CNN, DeepLabv3, Swin-UNet, and DINOv2, to assess their suitability for canopy segmentation under extreme data scarcity. Our experiments show that convolution-based models, particularly YOLOv11 and Mask R-CNN, generalize significantly better than transformer-based models. DeepLabV3, Swin-UNet and DINOv2 underperform likely due to differences between semantic and instance segmentation tasks, the high data requirements of Vision Transformers, and the lack of strong inductive biases. These findings confirm that transformer-based architectures struggle in low-data regimes without substantial pretraining or augmentation and that differences between semantic and instance segmentation further affect model performance. We provide a detailed analysis of training strategies, augmentation policies, and model behavior under the small-data constraint and demonstrate that lightweight CNN-based methods remain the most reliable for canopy detection on limited imagery.

1. Introduction

Tree canopy detection from aerial or satellite imagery plays a critical role in a wide range of environmental applications, including biomass estimation, carbon accounting, biodiversity monitoring, and urban forestry. Accurate canopy mapping enables more effective policy decisions and more scalable ecological assessment. Despite its importance, high-quality canopy annotations are expensive to acquire, and many real-world datasets consist of only a few hundred samples. This creates a challenging setting for modern deep learning models, which typically rely on large, diverse datasets to achieve robust performance.

The Solafune Tree Canopy Detection competition provides an example of this setting, offering only 150 labelled images for training. The dataset’s limited size, variability in lighting and terrain, and fine-grained boundaries of the canopy regions all contribute to the difficulty of the task. Additionally, individual trees are considered to be a distinct class from tight-knit groups of trees. This competition therefore, serves as an ideal testbed for examining how different model classes behave in extremely low-data environments.

In this paper, we investigate five modern deep learning approaches, YOLOv11, Mask R-CNN, DeepLabv3, Swin-UNet, and DINOv2, to evaluate their ability to localize and segment tree canopies. We focus on understanding why some architectures succeed while others fail in limited-data settings, highlighting the importance of inductive bias, pre-training, and model capacity when training from only 150 images. Our findings show that convolution-based models remain considerably more robust than transformer-based models without large-scale pretraining.

2. Background

Image Segmentation Semantic segmentation aims to assign a class label to every pixel in an image, making it well-suited for canopy mapping where fine-grained spatial boundaries must be localized. Classical approaches relied on hand-crafted features or region-based methods, but modern segmentation heavily depends on deep neural networks. Instance segmentation extends this by distinguishing between different instances of the same class within an image. It can be seen as a combination of object detection and semantic segmentation and is often approached with two-stage models, but computationally-efficient single-stage architectures have recently been the subject of increased advancement and popularity [17]. Two major approaches to segmentation tasks exist today: convolutional segmentation models that exploit spatial priors through localized filters, and transformer-based architectures that capture long-range context using global self-attention. In small-dataset settings such as Solafune, differences between these model classes become especially pronounced, motivating a comparative

study.

Remote Sensing Remote sensing imagery introduces unique challenges beyond standard computer vision tasks. Variations in spatial resolution, atmospheric conditions, sensor calibration, shadows, and land-cover diversity all contribute to high intra-class variability. Images may be acquired from satellites, fixed-wing aircraft, drones, or ground-based platforms, each providing different ground-sampling distances (GSD), focal lengths, and viewing geometries. As a result, the apparent size and shape of canopy features can vary dramatically across datasets, tree crowns may span only a few pixels in high-altitude satellite images, yet appear with fine structural detail in low-altitude UAV imagery. These factors create significant scale variation and require models that are robust to changes in feature size, texture, and resolution. Remote sensing datasets are also frequently limited in size due to the cost of expert annotation, making strong inductive biases and resistance to overfitting particularly important.

Convolution-based Approaches Convolutional Neural Networks (CNNs) such as U-Net, Mask-RCNN, DeepLabv3, and the YOLO family introduce useful spatial inductive biases such as locality and translation equivariance, making them naturally well-suited for segmentation when training data is scarce. Models like DeepLabv3 leverage dilated convolutions and multi-scale context through atrous spatial pyramid pooling (ASPP), enabling strong performance even with limited examples. Similarly, YOLOv11 offers a fast, anchor-free object detector with strong generalization due to extensive pretraining on large, diverse datasets. These properties help CNN-based models avoid severe overfitting on the 150-image Solafune dataset.

Transformer-based Approaches Vision Transformers (ViTs), including Swin-UNet and DINOv2, offer high modeling capacity and global receptive fields through self-attention mechanisms. However, ViTs lack the spatial priors inherent to CNNs and therefore require substantially more training data to learn stable visual representations. Swin-UNet introduces hierarchical attention and U-Net-like skip connections, but the core transformer blocks still remain highly data-dependent. DINOv2 provides powerful pretrained representations, yet fine-tuning with only 150 images leads to rapid overfitting and degraded segmentation quality. These limitations highlight the difficulty of applying transformer architectures to specialized remote sensing tasks without large-scale domain-relevant pretraining or heavy augmentation.

Tree Segmentation Traditional tree segmentation approaches rely heavily on LiDAR-derived 3D point cloud

data and are commonly divided into two categories. The first class consists of methods that convert the point cloud into a canopy height model (CHM) and then apply surface-based analysis to detect individual tree crowns, as demonstrated in Hui et al. [7] and Roussel et al. [12]. The second class comprises full 3D methods that operate directly on the raw point cloud. Among these, the point-based clustering algorithm, as proposed in Williams et al. [21] and Pang et al. [10], is a well-known and widely used method in this category due to its conceptual simplicity and its ability to capture tree-level structure.

The advent of deep learning has substantially expanded the methodological landscape for tree segmentation. For LiDAR-based workflows, Sun et al. [15] applies YOLOv4 and R-CNN architectures to perform crown segmentation on height maps derived from airborne LiDAR, while Wang et al. [20] uses R-CNN to identify and segment tree trunks from 3D point cloud representations. Beyond LiDAR, a growing body of work has explored segmentation directly from aerial imagery, often reporting improved scalability and generalization. High-resolution aerial data, such as RGB or SAR imagery, are particularly attractive because they are easier to acquire at scale and typically provide higher spatial resolution than LiDAR. For example, Velasquez-Camacho et al. [19] employs YOLOv5 and Fast R-CNN for tree segmentation on high-resolution aerial and satellite images, and Tolan et al. [18] introduces a self-supervised vision-transformer approach using DINOv2 to generate large-scale tree height maps from RGB imagery.

Although these methods demonstrate strong segmentation accuracy, most require extensive training datasets to achieve their reported performance. In contrast, there are works deployed on small datasets. For example, Takahashi et al. [16] highlights cases of ViT models outperforming CNNs in medical image analysis tasks with fewer required training images. Safanova et al. [13] discusses approaches for addressing small datasets in remote sensing, but does not investigate transformer-based architectures. However, to the best of our knowledge, no existing work explicitly investigates the suitability of modern deep learning architectures, including ViTs and CNNs, for remote sensing scenarios where only a small training dataset is available. This motivates our study, in which we evaluate which contemporary architectures are most effective under limited-data conditions.

3. Method

In this study, we focus on and compare the performance of five architectures: Yolov11 Seg, Mask R-CNN, DeepLabV3, Swin-UNet, and DINOv2. We randomly split the original training dataset into a training set and a validation set with a 4:1 ratio. To ensure reproducibility, we fix the random seed at 42.

Dataset The dataset provided by the competition [14] contains a training set and an evaluation (test) set, each with 150 aerial images in RGB tif format. Both datasets comprise five resolution groups - 10cm (38 images), 20cm (37 images), 40cm (25 images), 60cm (25 images), and 80cm (25 images). These images represent diverse land-cover scenes, including urban, rural, agricultural, industrial environments and open fields. All images have a uniform spatial resolution of 1024x1024 pixels. Samples of raw images and segmentation from the training dataset, with varying resolutions and scenes, are displayed in Figure 1. The competition evaluates solutions based on a weighted mean average precision (mAP) metric designed for instance segmentation tasks.

Yolo Ultralytics YOLO [8] offers a diverse suite of models designed to address a wide range of computer vision tasks, including object detection, instance segmentation, image classification, and pose estimation. YOLOv11 represents the latest iteration in the YOLO family, introducing architectural and training improvements aimed at enhancing both accuracy and efficiency. In this study, we begin our analysis with the lightweight YOLOv11-nano segmentation model (yolov11n-seg) and subsequently extend our experiments to the small, medium, and large variants to evaluate performance across model scales. All models are initialized from weights pretrained on the COCO dataset, providing a strong general-purpose feature representation. We then fine-tune these pretrained models on the custom canopy tree segmentation dataset to adapt them to the domain-specific characteristics of aerial forest imagery. In addition, the training configuration uses an input size of 1024x1024, allowing images to be fed into the network without any resizing. All YOLOv11 segmentation models are trained with the AdamW optimizer with a learning rate of 1e-4 for a total of 100 epochs.

Mask R-CNN Mask R-CNN [6] is used as a baseline fully-convolutional network. This model was designed specifically for instance segmentation and has been successfully applied across domains [2]. Unlike YOLO, it is a two-stage detector consisting of a region proposal network followed by classification. Weights are initialized from the pretrained model as obtained in the original paper. The pre-training utilized the COCO dataset with images rescaled to several resolutions ranging from 640 to 800 pixels in width and height. The model is fine-tuned on the tree canopy dataset with random cropping and flipping augmentations. Additionally, the input images are resized to a height and width of 640 pixels to match the pretraining data and reduce computational requirements. A learning rate of 1e-4 is used with the Adam optimizer.

DeeplabV3 DeepLabV3 [4] was used to gauge the performance of semantic segmentation with a convolution-based architecture. This model, developed the same year as Mask R-CNN, innovated by incorporating atrous convolutions and atrous spatial pyramid pooling into a fully convolutional network. These additions enable scale-invariant learning, while avoiding computational costs associated with high-resolution training. The final classification layer was adjusted to output logits for three classes: background, individual tree, and group of trees. In order to apply instance segmentation metrics, the output semantic segmentation masks are converted to polygons using the OpenCV Python library. We load weights from the model pretrained on ImageNet at a 224x224 resolution. The input images are resized to match the pretraining resolution, and randomly augmented with Gaussian blur and noise, brightness and contrast changes, flips, and rotations. The Adam optimizer is used with a learning rate of 1e-4.

DinoV2 For the transformer-based baseline, we implement a lightweight semantic segmentation model using the DINOv2-Base vision transformer as a frozen encoder. Because DINOv2 does not include a built-in decoder for dense prediction, we adapt the model by attaching a 1×1 convolutional segmentation head to the final patch embeddings, followed by bilinear upsampling to produce class logits at a 224×224 resolution. Only this segmentation head is trained; the 86M-parameter DINOv2 backbone remains frozen to avoid overfitting and to stay within the memory constraints of the GPU.

All training and evaluation images are resized to 224×224 and normalized using the ImageNet mean and standard deviation expected by DINOv2. During training, we apply standard augmentations, including flips, rotations, Gaussian blur and noise, to counteract the extreme data scarcity of the 150 image dataset. The model is optimized using Adam optimizer with a learning rate of 1e-4, cross-entropy loss, and training is run for 100 epochs.

Swin-UNet For a hierarchical transformer-based encoder-decoder model, we use Swin-UNet, following the original architecture and official implementation provided by Cao et al. [3]. Swin-UNet extends the Swin Transformer to dense prediction tasks using a symmetric U-shaped design with skip connections and shifted window self-attention, allowing the model to capture both local and global context efficiently. In our setup, we adopt the Swin-UNet-Tiny configuration and initialize the backbone with ImageNet-pretrained Swin Transformer weights. The encoder processes features through four hierarchical stages, while the decoder reconstructs spatial resolution through patch-expansion layers and skip connections from corresponding encoder stages.



Figure 1. Examples of Training Images with Different Resolutions and Scenes and their Corresponding Segmentation

All images are resized to 224x224 to match the model’s native patch size (4x4). We apply standard augmentations, including random flips and rotations, to improve robustness. The model is trained using the AdamW optimizer with a learning rate of 1e-4 and cross-entropy loss for 150 epochs with a batch size of 16. This configuration enables Swin-UNet to effectively capture fine canopy boundaries and small-scale structures in aerial forest imagery.

4. Results

The weighted mAP scores obtained on the Solafune Tree Canopy Detection validation (test) set from the five experimental models are reported in Table 1. As shown in the table, YOLOv11 and Mask R-CNN substantially outperform the other three architectures, DeepLabv3, Swin UNet and DINOv2, demonstrating a clear advantage in segmentation accuracy. This performance gap is further illustrated in Figure 12, which presents qualitative comparisons across five representative examples from the evaluation set and highlights the superior delineation quality achieved by YOLOv11 and Mask R-CNN.

As illustrated in the figure, both YOLOv11 and Mask R-CNN exhibit consistently strong performance across the five representative examples, indicating robust generalization across diverse scene types. In contrast, DeepLabv3, Swin UNet and DINOv2 show noticeably weaker performance, particularly in areas containing densely clustered trees or agricultural landscapes, while tending to perform comparatively better on images from urban environments. This suggests that these architectures may struggle with complex canopy structures and exhibit a bias toward more homogeneous or structurally regular scenes. The detailed analysis for each model is as follows:

Table 1. Weighted mAP (IoU-based) on the Solafune Tree Canopy Detection test set.

| Model | Weighted mAP |
|-------------------|--------------|
| YOLOv11 Seg Large | 0.281 |
| Mask R-CNN | 0.219 |
| DeepLabv3 | 0.038 |
| Swin-UNet | 0.022 |
| DINOv2 | 0.021 |

Yolo The performance for four Yolov11 segmentation models is shown in Figure 2 and Figure 3. The small model exhibits early stopping at approximately 80 epochs upon reaching its optimal validation performance, likely due to its moderate capacity, which enables it to capture the key features in the dataset more efficiently than the Nano model. In contrast, the other models continue training for the full 100 epochs, as the Nano model converges more slowly due to limited capacity, and the larger models require more epochs to fully leverage their increased representational capacity. Overall, the experiments show a clear trend of improved mask prediction performance with increasing model capacity. Training and validation losses decreased steadily across all runs, and larger models generally achieved higher mask mAP.

A noticeable gap exists between validation and test mAP, reported in Table 2, with test performance consistently higher despite similar or slightly lower validation mAP. This discrepancy is likely due to the relative size of the datasets: the validation set is much smaller than the test set, making its mAP estimates more variable and sensitive to the specific samples it contains. With only 30 images in validation versus 150 in testing, the validation mAP may not fully capture the model’s generalization ability, whereas the

Table 2. Weighted (Test) mAP (IoU-based) for YOLOv11 Models

| Model | Val mAP | Test mAP |
|--------------------|---------|----------|
| YOLOv11 Seg Large | 0.189 | 0.281 |
| YOLOv11 Seg Medium | 0.187 | 0.279 |
| YOLOv11 Seg Small | 0.177 | 0.257 |
| YOLOv11 Seg Nano | 0.164 | 0.249 |

larger test set provides a more stable and representative assessment. Additionally, the competition uses a customized implementation of mAP that applies different weightings to certain scenes. Among the four models, the largest variant achieved the highest test mAP, indicating that increased model capacity better captures the complexity of the segmentation task, even if validation performance alone might not fully reflect this advantage.

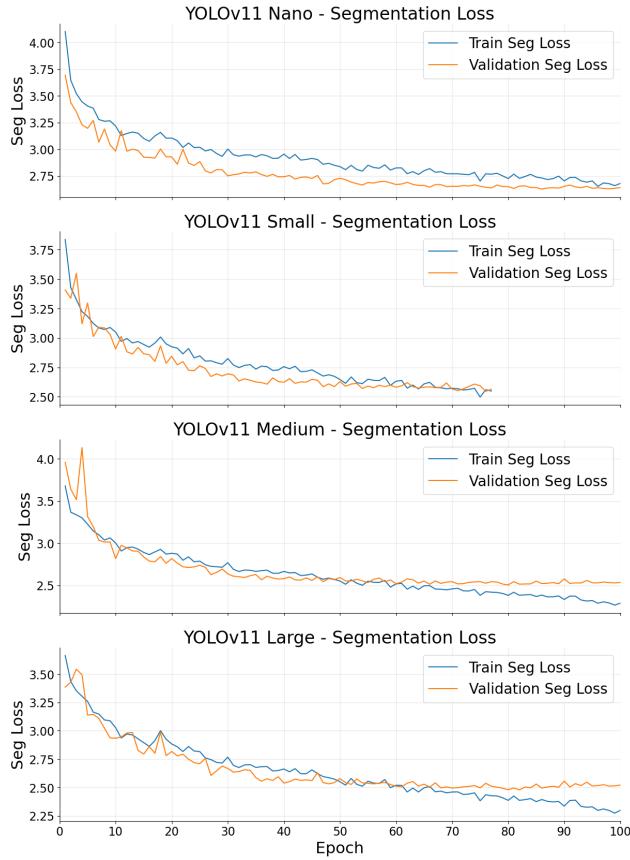


Figure 2. Training and validation loss curves for the YOLOv11 Seg models

Mask R-CNN Figure 4 and Figure 5 show the performance of Mask R-CNN during training. The loss is stable, smoothly increasing with a plateau at about 15 epochs.

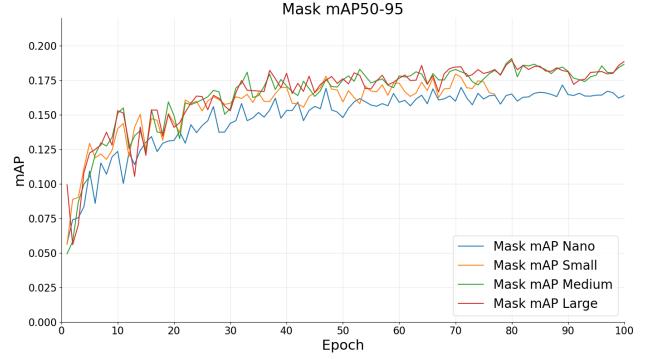


Figure 3. Validation (Mask) mAP for the YOLOv11 Seg Models

The validation mAP further emphasizes the early plateau. This may be a sign of a difficult to traverse region of the loss landscape, and could indicate potential for improvement with further hyper-parameter tuning. The training and validation performance translated well to the test dataset, achieving a reasonable weighted mAP of 0.22. Similarly to Yolo, this result is higher than the final mAP of 0.14 achieved on the validation set after training.

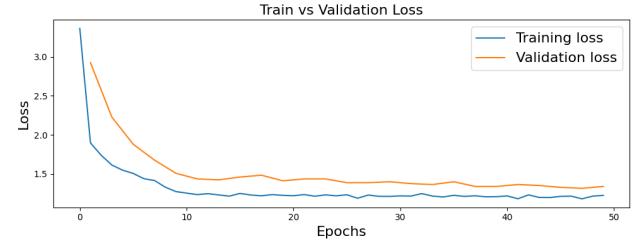


Figure 4. Training and validation loss curves for the Mask R-CNN model

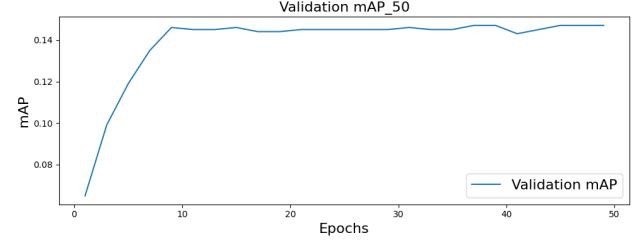


Figure 5. Validation mAP curve for the Mask R-CNN model

DeeplabV3 The DeeplabV3 model steadily improved throughout training, as shown by the validation loss in Figure 6 and pixel accuracy in Figure 7. The number of epochs observed was increased from 100 to 300 epochs as validation loss indicated further potential gains without overfitting. The model achieved a final pixel accuracy of 0.82 and

mean IoU of 0.71 on the validation set. This indicates successful learning and generalization for semantic segmentation. However, running inference on the test set and converting the masks to polygons resulted in a low weighted mAP score of 0.038.

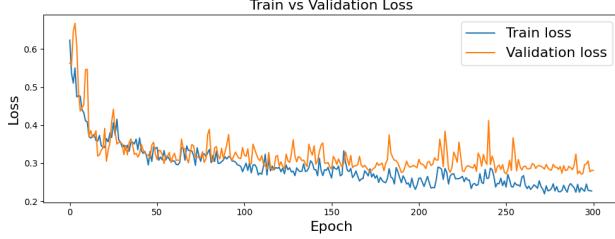


Figure 6. Training and validation loss curves for the DeepLabV3 model

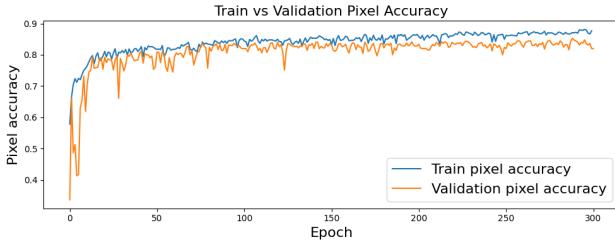


Figure 7. Training and validation pixel accuracy curves for the DeepLabV3 model

DinoV2 As seen in Figure 9 and Figure 8, the DINOv2 segmentation head trains stably, with loss decreasing and pixel accuracy improving consistently over the training epochs. This indicates that the model was able to learn a meaningful mapping on the training split despite relying on a frozen backbone and a very small dataset. However, when applied to the evaluation set, the predictions failed to transfer effectively. The final weighted mAP remained extremely low, largely because the evaluation images exhibited distributional differences and greater variability than the 150 training examples could represent. As a result, DINOv2 produced sparse or incomplete canopy predictions that translated poorly into polygon-based instance masks. These findings suggest that, unlike CNN-based methods, transformer models such as DINOv2 require substantially more data or stronger domain-specific pretraining to perform well in small, heterogeneous remote-sensing datasets.

Swin-UNet To evaluate transformer-based segmentation performance on the tree canopy detection task, we trained a Swin-UNet model for 150 epochs using the same training/validation split and preprocessing strategy applied to

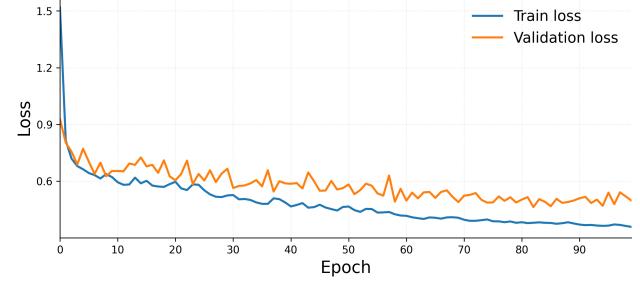


Figure 8. Training and validation loss curves for the DINOv2 model

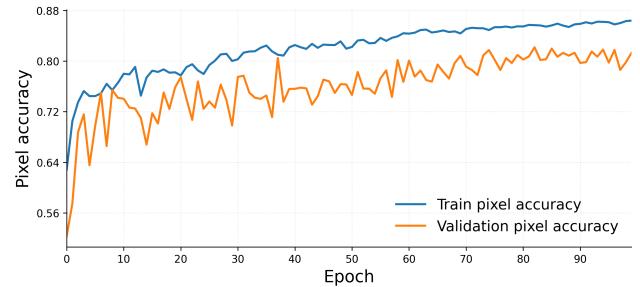


Figure 9. Training and validation pixel accuracy curves for the DINOv2 model.

the other baselines. The training process was stable, with both training and validation losses decreasing smoothly over time, as shown in Figure 10. Despite the relatively small dataset size, the model converged without oscillation or divergence, indicating that the hybrid CNN-Transformer architecture was able to learn meaningful spatial features.

Pixel accuracy curves (Fig. 11) show that the model achieved steady improvements throughout training: training pixel accuracy increased to roughly 0.87, while validation pixel accuracy stabilized around 0.85. This suggests that Swin-UNet generalized reasonably well to the held-out validation set in terms of pixel-wise correctness.

However, these promising accuracy and loss trends did not translate into strong instance-level performance on the test set. As shown in Table 1, Swin-UNet produced a very low weighted mAP. The primary reason is that mAP is computed from polygonized instance masks, which penalizes even small inconsistencies at object boundaries. Although the model captured broad canopy regions, the predicted masks tended to be fragmented or overly smooth, causing large errors once converted to polygon instances. This highlights a key limitation of transformer-based architectures in this setting: they require larger, more diverse datasets or domain-specific pretraining to produce crisp and topologically consistent segmentation masks suitable for polygon extraction.

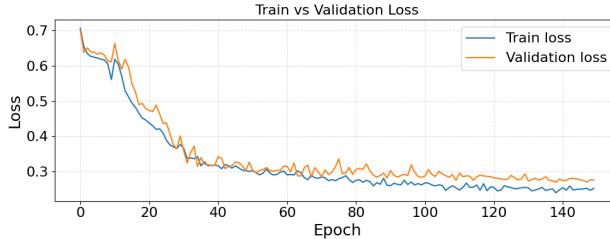


Figure 10. Training and validation loss curves for the Swin-UNet model.

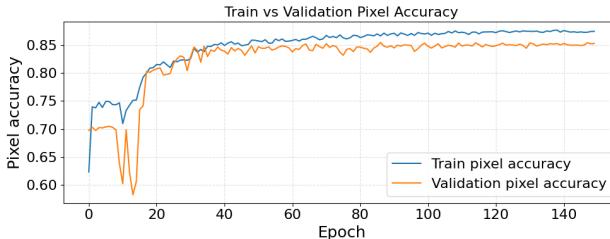


Figure 11. Training and validation pixel accuracy curves for the Swin-UNet model.

5. Discussion

DinoV2 A key reason for DINOv2’s poor performance in our setting is the extremely limited amount of training data available. The Solafune dataset contains only 150 annotated images, which is far too small for a high-capacity Vision Transformer to learn meaningful canopy representations. This contrasts sharply with recent work demonstrating DINOv2’s success in tree-mapping tasks, such as Tolan et al. [18], where the model was pretrained on 18 million unlabeled aerial images and then refined using 15 thousand of domain-specific canopy masks. In that setting, both the scale of the data and the close alignment between the pre-training imagery and the target task were critical to achieving strong segmentation performance.

In our experiments, neither of these conditions held: the dataset was two orders of magnitude smaller, and the evaluation imagery differed noticeably from the training set in appearance, resolution, and acquisition conditions. As a result, the frozen DINOv2 backbone lacked both sufficient data and domain context to transfer effectively. We therefore expect that, had a substantially larger and more homogeneous set of canopy images been available, similar in scope and structure to the datasets used in prior DINOv2 remote-sensing studies, the transformer would have produced far more accurate and stable canopy predictions on the Solafune competition.

Swin-UNet Beyond DINOv2 specifically, our results show a substantial performance gap between Vision Trans-

former-based models and conventional CNN architectures. While YOLOv11 segmentation models achieved weighted mAP scores between 0.25 and 0.28, both Swin-UNet and DINOv2 were near 0.02. This gap is consistent with well-established findings that Vision Transformers require significantly more data to achieve competitive performance, whereas CNNs are far more sample-efficient due to their strong inductive biases [5, 11].

CNNs incorporate spatial priors such as locality, translation equivariance, and hierarchical feature extraction, enabling them to generalize reliably even with small or heterogeneous datasets. In contrast, ViTs rely on self-attention without built-in spatial priors and therefore depend heavily on large-scale, domain-aligned pretraining to learn low-level structure [11]. Without such pretraining—and with only 150 labeled images—the transformer models in our study underfit severely.

This explains why Swin-UNet did not perform well despite its hierarchical attention design [9]. The model struggled to learn stable canopy boundaries under data scarcity and domain shift, producing noisy segmentation outputs. Meanwhile, CNN-based models were better able to extract texture, edge, and shape cues from limited imagery, leading to substantially higher accuracy. Overall, our findings align with prior work showing that ViTs become competitive in remote sensing only when supported by large-scale or domain-specific pretraining [1].

Significance of Architectural Differences The two main properties varying across the examined architectures are the choice between a CNN or ViT design, and the specific type of segmentation head. In general, CNNs outperformed ViT models, which we attribute to ViT’s limited ability to learn from small datasets. Another significant aspect of the models was the segmentation head, which varied between designs for instance and semantic segmentation. During training, the semantic models appeared to learn effectively, evidenced by stable loss curves and validation metrics such pixel accuracy. However, testing on the evaluation set revealed poor performance on the weighted mAP metric designed for instance segmentation tasks.

The disparity between architectures designed for different tasks is best demonstrated by comparing Mask R-CNN with DeeplabV3. Both models were introduced in the same year, achieving state-of-the-art performance on the instance and semantic segmentation tasks for which they were designed. However, our results show a severe disparity between the two according to the Weighted mAP metric.

It is clear that semantic segmentation models do not generalize well to the instance segmentation task. This may bias the comparison of CNN and ViT models, as neither of the tested ViT architectures were originally designed with an instance segmentation head. Thus, to further support our

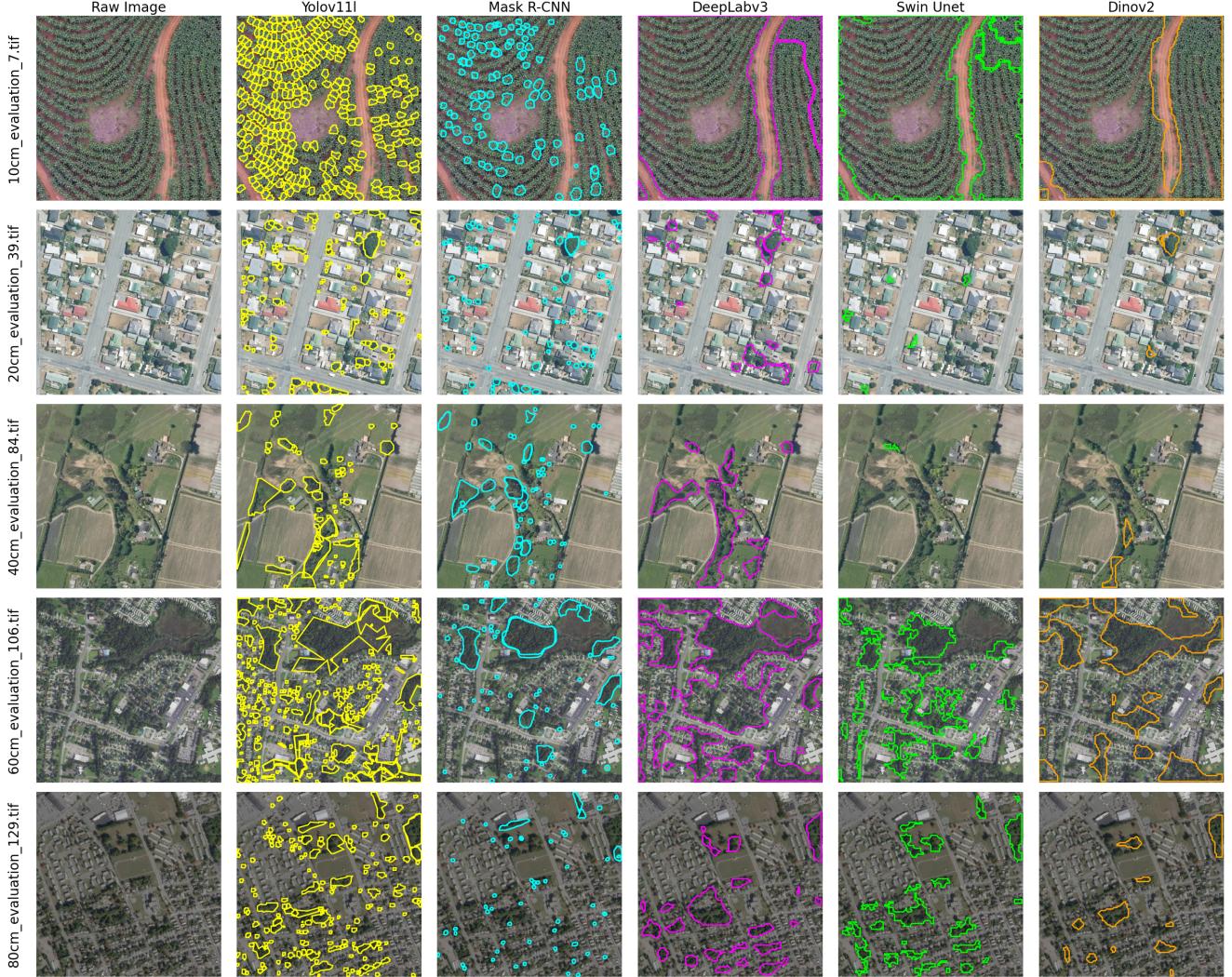


Figure 12. Segmentation Visualization for the Five Architectures

analysis of ViTs and CNNs, we highlight the pixel accuracies achieved by DeepLabV3 and DinoV2. At 100 epochs, DeepLabV3’s accuracy is higher than that of DinoV2 and continues to grow while DinoV2 begins to overfit. Even when comparing a modern ViT architecture to an older CNN with a metric not specific to either kind of segmentation, the CNN is observed to learn better from the small dataset.

6. Conclusion

We evaluated five deep learning architectures on a small-scale remote-sensing dataset and compared their effectiveness under limited training data. Among the models, the convolution-based approaches, YOLOv11 and Mask R-CNN, achieved the most stable and accurate performance on the 150-image Solafune dataset. Their inherent spa-

tial inductive biases, together with extensive pretraining on large-scale image corpora, enabled them to generalize effectively despite scarce supervision.

In contrast, both Vision Transformer approaches, SwinUNet and DINOV2, exhibited pronounced overfitting and consistently failed to produce reliable canopy delineations. DeepLabv3 showed similar weak performance, likely due to the combination of its semantic-segmentation formulation and the limited size of the training set. Both factors hinder its ability to capture fine canopy boundaries.

Overall, the results indicate that in extremely low-data remote-sensing settings, lightweight CNN architectures and instance-segmentation frameworks remain considerably more robust than ViT-based or purely semantic segmentation designs. Transformer models appear to require substantially larger datasets or domain-specific pre-training to achieve competitive performance in this setting.

References

- [1] Chamira Bandara, Vishal Patel, et al. Transformers in remote sensing: A survey. *Remote Sensing*, 14(15):3366, 2022. 7
- [2] Sreya Ramesh C and Vinod Kumar V. A review on instance segmentation using mask-rcnn. *SSRN Electronic Journal*, 2020. 3
- [3] Hu Cao, Yue Wang, Jiarui Chen, Dongsheng Jiang, Xin Zhang, Qi Tian, and Yunhai Wang. Swin-unet: Unet-like pure transformer for medical image segmentation. *arXiv preprint arXiv:2105.05537*, 2021. 3
- [4] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation, 2017. 3
- [5] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, and et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. 7
- [6] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn, 2018. 3
- [7] Zhenyang Hui, Penggen Cheng, Bisheng Yang, and Guoqing Zhou. Multi-level self-adaptive individual tree detection for coniferous forest using airborne lidar. *International Journal of Applied Earth Observation and Geoinformation*, 114: 103028, 2022. 2
- [8] Glenn Jocher and Jing Qiu. Ultralytics yolo11, 2024. 3
- [9] Ze Liu, Yutong Lin, Yue Cao, and et al. Swin transformer: Hierarchical vision transformer using shifted windows. In *IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021. 7
- [10] Yong Pang, Weiwei Wang, Liming Du, Zhongjun Zhang, Xiaojun Liang, Yongning Li, and Zuyuan Wang. Nyström-based spectral clustering using airborne lidar point cloud data for individual tree segmentation. *International Journal of Digital Earth*, 14(10):1452–1476, 2021. 2
- [11] Aravindh Raghu, Thomas Unterthiner, Simon Kornblith, and et al. Do vision transformers see like convolutional neural networks? In *Advances in Neural Information Processing Systems*, 2021. 7
- [12] Jean-Romain Roussel, David Auty, Nicholas C. Coops, Piotr Tompalski, Tristan R.H. Goodbody, Andrew Sánchez Meador, Jean-François Bourdon, Florian de Boissieu, and Alexis Achim. lidr: An r package for analysis of airborne laser scanning (als) data. *Remote Sensing of Environment*, 251:112061, 2020. 2
- [13] Anastasiia Safonova, Gohar Ghazaryan, Stefan Stiller, Magdalena Main-Knorn, Claas Nendel, and Masahiro Ryo. Ten deep learning techniques to address small data problems with remote sensing. *International Journal of Applied Earth Observation and Geoinformation*, 125:103569, 2023. 2
- [14] Solafune, Inc. Tree canopy detection — competition overview. <https://solafune.com/competitions/26ff758c-7422-4cd1-bfe0-daecfc40db70?menu=about&tab=overview>, 2025. Accessed: 2025-10-09. 1, 3
- [15] Chenxin Sun, Chengwei Huang, Huaiqing Zhang, Bangqian Chen, Feng An, Liwen Wang, and Ting Yun. Individual tree crown segmentation and crown width extraction from a heightmap derived from aerial laser scanning data using a deep learning framework. *Frontiers in Plant Science*, Volume 13 - 2022, 2022. 2
- [16] Satoshi Takahashi, Yusuke Sakaguchi, Nobuji Kouno, Ken Takasawa, Kenichi Ishizu, Yu Akagi, Rina Aoyama, Naoki Teraya, Norio Shinkai, Hidenori Machino, Kazuma Kobayashi, Ken Asada, Masaaki Komatsu, Syuzo Kaneko, Masashi Sugiyama, and Ryuji Hamamoto. Comparison of vision transformers and convolutional neural networks in medical image analysis: A systematic review. *Journal of Medical Systems*, 48:84, 2024. 2
- [17] Di Tian, Yi Han, Biyao Wang, Tian Guan, Hengzhi Gu, and Wei Wei. Review of object instance segmentation based on deep learning. *Journal of Electronic Imaging*, 31, 2021. 1
- [18] Jamie Tolan, Hung-I Yang, Benjamin Nosarzewski, Guillaume Couairon, Huy V. Vo, John Brandt, Justine Spore, Sayantan Majumdar, Daniel Haziza, Janaki Vamaraju, Theo Moutakanni, Piotr Bojanowski, Tracy Johns, Brian White, Tobias Tiecke, and Camille Couprie. Very high resolution canopy height maps from rgb imagery using self-supervised vision transformer and convolutional decoder trained on aerial lidar. *Remote Sensing of Environment*, 300:113888, 2024. 2, 7
- [19] Luisa Velasquez-Camacho, Maddi Etxegarai, and Sergio de Miguel. Implementing deep learning algorithms for urban tree detection and geolocation with high-resolution aerial, satellite, and ground-level images. *Computers, Environment and Urban Systems*, 105:102025, 2023. 2
- [20] Jiamin Wang, Xinxin Chen, Lin Cao, Feng An, Bangqian Chen, Lianfeng Xue, and Ting Yun. Individual rubber tree segmentation based on ground-based lidar data and faster r-cnn of deep learning. *Forests*, 10(9), 2019. 2
- [21] Jonathan Williams, Carola-Bibiane Schönlieb, Tom Swinfield, Juheon Lee, Xiaohao Cai, Lan Qie, and David A. Coomes. 3d segmentation of trees through a flexible multi-class graph cut algorithm. *IEEE Transactions on Geoscience and Remote Sensing*, 58(2):754–776, 2020. 2