

# Tree Instance Segmentation with Temporal Contour Graph

Adnan Firoze Cameron Wingren<sup>†</sup> Raymond A. Yeh Bedrich Benes Daniel Aliaga  
 Department of Computer Science, Purdue University  
 Department of Forestry and Natural Resources<sup>†</sup>, Purdue University  
 {afiroze, cwingren, rayyeh, bbenes, aliaga}@purdue.edu

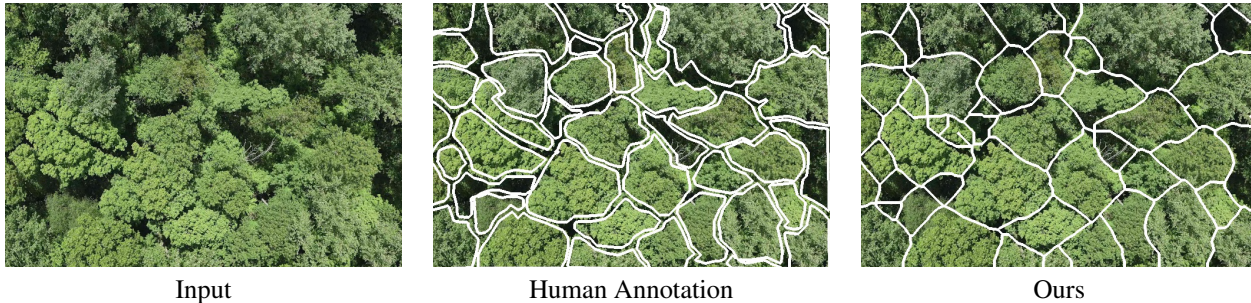


Figure 1. **Tree Instance Segmentation.** A flyover from a UAV and the corresponding tree instance segmentation during green leaf season when the tree crowns are most similar to their neighbors, and occlusion is complex adding to the difficulty in segmentation.

## Abstract

We present a novel approach to perform instance segmentation and counting for densely packed self-similar trees using a top-view RGB image sequence. We propose a solution that leverages pixel content, shape, and self-occlusion. First, we perform an initial over-segmentation of the image sequence and aggregate structural characteristics into a contour graph with temporal information incorporated. Second, using a graph convolutional network and its inherent local messaging passing abilities, we merge adjacent tree crown patches into a final set of tree crowns. Per various studies and comparisons, our method is superior to all prior methods and results in high-accuracy instance segmentation and counting despite the trees being tightly packed. Finally, we provide various forest image sequence datasets suitable for subsequent benchmarking and evaluation captured at different altitudes and leaf conditions.

## 1. Introduction

Trees in forests are tightly spaced, partially overlapping 3D objects with complex boundaries. Tree instance segmentation is critical in several domains. For example, ecosystem services and agriculture need to segment and count trees in large areas in order to obtain information about the ecological balance, environmental health, and timber inventory. Counting trees from the ground per-

spective is inefficient, does not scale, and is challenging to automate because of many occlusions with branches and low accessibility. In this paper, we address tree instance segmentation and counting using overhead RGB image sequences captured by unmanned aerial vehicles (UAVs), especially during the green-leaf season when trees are most self-similar; see Fig. 1 for an illustration.

There is significant prior work in segmentation and counting, particularly in the field of instance segmentation. While some approaches make use of LiDAR or RGB-D images (see the survey paper [5]), we focus on using easier-to-obtain uncalibrated RGB image sequences. Prior works based on uncalibrated RGB images can be largely organized into three groups. The first group seeks to count individual objects. These approaches often use density estimation and do not focus on segmentation [35, 36]. The second group of methods relies on convolutional neural networks (CNNs) applied directly to image pixels, such as Mask R-CNN [26]. However, in the case of abutting and self-similar objects, e.g., trees, distinguishing individual instances is hard. The third group of techniques (e.g., Ke et al. [28], Newell and Deng [43]) model object contour as a graph, where each pixel corresponds to a node, and makes use of graph convolutional networks to complete individual contours; nevertheless, abutting and self-similar instances also hinder these methods. Yet other methods, such as those in digital forestry research, exploit domain-specific features such as assumed differences between tree species or fall leaf coloring (i.e., during one brief time period of the year, the

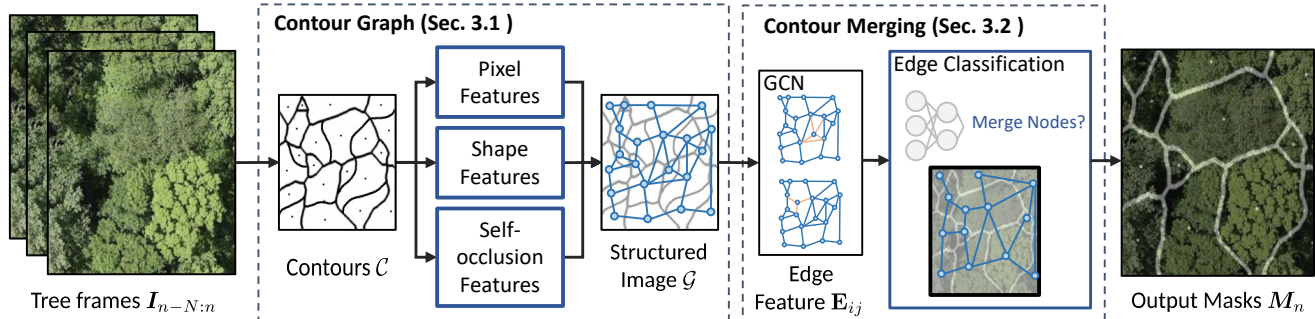


Figure 2. **Pipeline:** The input image sequence is analyzed. Initial contours and features are detected and organized into a contour graph that is refined by merging edges and nodes, resulting in the final output mask.

leaf color of adjacent crowns is often different).

Our approach is motivated by a key observation that two tree crown leaf patterns tightly packed together are highly similar, and additional features beyond leaf patterns are necessary to perform segmentation. Hence, we consider features based on the tree crown shape because it is unlikely to observe a rectangular tree crown. Beyond shape features, we also use the changing self-occlusion patterns captured in the different frames to aid segmentation.

At a high level, our proposed approach simultaneously exploits pixel content, shape, and self-occlusion, which collectively define a graph-based structure that we call a *contour graph*. Each node corresponds to an initial over-segmentation of a tree crown fragment; i.e., each node corresponds to a *region enclosed* by a closed contour. We then learn features on this contour graph via graph convolutional network (GCN) to determine which nodes should be merged. Our method decides whether two nearby regions correspond to the same tree crown. Notably, a tree crown fragment is subject to various simultaneous features that we can exploit to discern one tree crown from another, even if one tree crown is of the same species and has a very similar leaf pattern and color to an adjacent tree crown. Altogether this leads to an instance segmentation method that can process overhead RGB image sequences of dense forests even when all leaves are mostly similar in color during summer. See Fig. 2 for an overview of our approach.

When creating and evaluating model performance, we are not aware of suitable databases on dense forests with subsequent frames. To address this: (a) We leveraged developmental tree models [33, 57] to produce a synthetic dataset with annotated tree crowns (5,157 trees in total); (b) We manually labeled real-world image sequences captured by UAV over three large forests (6,527 trees in total), collectively spanning approximately 3,680,000 m<sup>2</sup>. We will make our self-collected datasets publicly available.

On these datasets, we show that the proposed method achieves a segmentation accuracy of 73.6 and a count accuracy of 89.8% on average, which is compared to multiple recent instance segmentation approaches.

### Our main contributions include:

- an instance segmentation method to robustly process densely packed trees where the instances are abutting, partially overlapping, and self-similar,
- tree crown counting, which is beneficial to ecosystem services and digital forestry, and
- a curated dataset of multiple labeled and unlabeled temporally continuous dense forests suitable for future research.

## 2. Related Work

**Counting Methods.** Instance segmentation and counting have been pursued by various approaches, including methods using input data beyond uncalibrated RGB image sequences. Counting methods use stochastic mechanisms to estimate the number of instances. For example, Lian et al. [35] use RGB-D data and density estimation to distinguish instances at a relatively low resolution of individuals in a dense crowd of humans. Firoze et al. [18] perform a density-based estimation of trees but use 12 months of satellite data. Liu et al. [36] use RGB-thermal imagery to perform crowd counting using a multi-modality deep network and exploiting the thermal signatures of humans. In our case, each instance is relatively larger (i.e., occupies many pixels). Thus resolving them is not a challenge because of limited resolution but rather difficult due to similarity, semi-transparency, tight spacing, and partial overlaps with complex boundaries. A density-based method would be too inaccurate and would not provide a segmentation.

**Detection and Instance Segmentation.** Object detection methods largely fall into two-stage and single-stage methods. Two-stage detection, such as Faster R-CNN [51] and Fast R-CNN [19], predict object masks based on region proposal and bounding box regression heads using features extracted from convolutional neural networks (CNN). Single-stage detection [31, 37, 50, 55, 56], bypasses the regional proposal stage. Notably, YOLO-based methods [9, 34, 50, 58] achieve efficient detection of objects in real-time. Nevertheless, when applied to abutting and self-similar content,

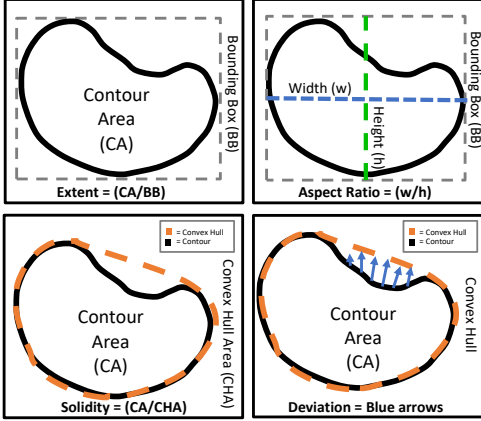


Figure 3. **Shape Features.** Illustration of the extracted shape features, including extent, aspect ratio, solidity, and deviation.

e.g., tree crowns, the precision of these approaches is low for tree counting.

Many instance segmentation methods [10, 11, 20, 25, 26, 30, 32, 62] have been developed. Earlier works use CNN features and build on the detection methods, e.g., Mask R-CNN [26]. More recently, several transformer-based instance segmentation methods [13, 40, 60] have been proposed as well, for instance, for instance, SwinTransformer [39, 40] uses a transformer-based backbone together with Mask R-CNN. Some methods use graph convolutional networks (GCN) to perform instance segmentation [28, 32]. For example, Ke et al. [28] models a graph corresponding to a sequence of pixels (nodes) defining the contour (connectivity) of object instances aiming to model occlusions among the objects. Different from our approach, we define a node to be the “region” enclosed by a contour and exploit the change in “self-occlusion”, through time, within each object (i.e., tree crown). In this work, we compare several of the recent instance segmentation methods, including Mask-RCNN with ResNet and Swin-T backbone [26, 40], TraDes [59], and BoundaryFormer [32].

**Digital Forestry and Remote Sensing** often use CNNs (e.g., see summary by Chen et al. [12] on remote sensing) and exploit domain-specific characteristics. One methodology is to acquire LIDAR data from “underneath” the tree crowns using UAVs [2]. The data can be used to measure trunk diameter and counts, but unfortunately, it does not scale as flying through large dense forests is a significant challenge. Another strategy is to use NDVI, an image-based vegetation index, to distinguish vegetation from non-vegetation. However, this index cannot separate one tree crown from another. Another application is tree classification, which also relies on CNNs [45] and the instance segmentation of trees is suggested as future work.

Other approaches exploit the visual differences between tree species. For example, Liu et al. [38] uses a CNN and UAV captured overhead RGB data to segment and classify

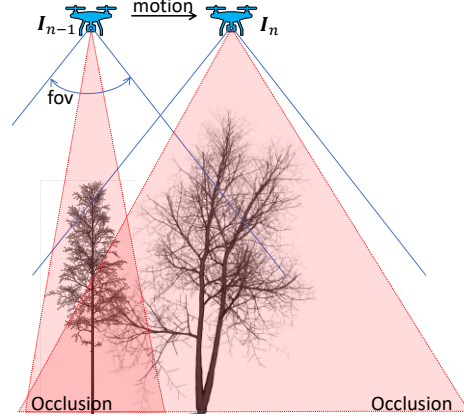


Figure 4. **Self-Occlusion.** Consecutive frames  $I_{n-1}$  and  $I_n$  from the moving UAV capture slightly different self-occlusion patterns.

trees. However, the segmentation component of this work is done in the peak fall season (when leaves of adjacent trees are most likely a different color) and then subject to *manual correction* to determine a good set of segmented trees. They claim accurate instance segmentation of trees as a challenging task for future work.

**Tree modeling** has a long history in computer vision and computer graphics [49]. Early methods focused on fractal-based approaches [1, 4, 52], and later methods simulate botanical plant model [16] development [6, 8, 21], including competition for resources [42, 46], climbing vegetation [24], ecosystems [44], or plants interacting with wind [47] or fire [23]. Synthetic tree models are perceived as highly realistic [48] and provide sufficient details to bridge the simulation-to-real gap for vision-based approaches. We leverage these tree models to help develop the approach presented in this paper.

### 3. Tree Crown Instance Segmentation

We formulate our problem as computing a set of instance segmentation masks for an uncalibrated input image sequence. The input image in a sequence is denoted by  $I_k \in (I_1, \dots, I_N)$  and the corresponding set of predicted instance segmentation masks for  $I_k$  is denoted by  $M_k = \{M_{t,k} \forall t \in \mathcal{T}_k\}$  for the tree crowns  $t \in \mathcal{T}_k$  in the current frame  $I_k$ . Each segmentation mask of a tree crown is represented by a set of pixels  $M_{t,k} = \{(u, v)\}$ . Our approach creates a contour graph representation capturing an initial over-segmentation of the forest and formulates the prediction of tree crown masks as merging the over-segmented regions.

#### 3.1. Contour Graph Creation

We create a graph-based data structure  $\mathcal{G} = (\mathcal{C}, \mathcal{E})$  (*contour graph*) for assisting with differentiating tightly-spaced and highly-similar tree crowns. The contour graph readily

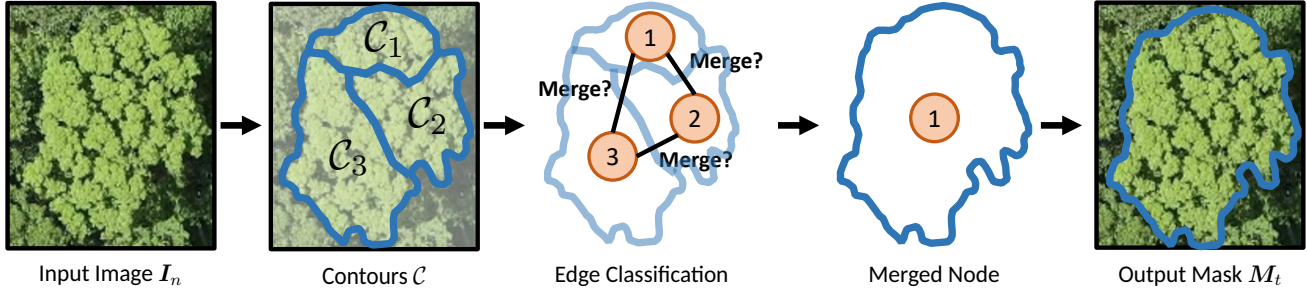


Figure 5. **Merging Process:** Illustrative presentation of our process to/from contour and graph spaces to merge noisy contours using edge classification (detailed in Sec. 3.2).

stores different types of feature sets and supports merging nodes via edge collapse operations.

**Contour Graph Construction.** We create a graph where each node corresponds to a closed contour  $C_t = \{(u, v)\}$ , i.e., the set of pixels of a tree crown fragment in the image. Two contours (nodes) have an edge between them  $(C_i, C_j) \in \mathcal{E}$ , if they are adjacent to each other in the image.

We extract the initial contours by using deep edge detection [53] that produces an edge map with plausible tree crown fragments. The extracted edge map is often noisy so we apply Guo-Hall algorithm [22] to skeletonize the edge map. Then, a simple polygonal contour can be used to represent the pixels of each closed region (i.e., each tree crown fragment). In the following, we discuss how to extract three feature set types that are stored within each node.

**Pixel Features** represent the pixel content within a contour. We use the following features: *Patch pixels*: we extract the center  $p \times p$  RGB image patch within each contour, using  $p = 30$ ; *Pixel similarity*: for a contour, we compute the similarity of its center patch to the center patches of neighboring contours. We measure similarity using LPIPS [61] which corresponds highly to human perception.

**Shape Features.** We capture the geometric shape features of a contour (see Fig. 3) including its *area*, i.e.,  $|C_i|$ , *extent* given as the ratio of contour area to its bounding rectangle area, *aspect ratio* of  $C_i$ 's bounding box, *solidity* given by the ratio of the area to the contour's convex hull, and *deviation* of a contour from its convex hull.

**Self-occlusion Features.** We capture the self-occlusion pattern of a contour. Fig. 4 depicts a scenario of two adjacent trees. The pattern of how the leaves and branches of the left tree exhibit their occlusion relationship as the image sequence was captured is likely different than the pattern of the right tree. This is true for trees of different species but also for trees of similar species. Although both branching patterns of adjacent trees originate from approximately the middle of the tree crown, they are not the same pattern. The optical flow of the pixels within each tree crown fragment can be compared and thus used to help differentiate one self-occlusion pattern from another.

We use the Gunnar-Farneback [17] optical flow algorithm between images  $I_n$  and  $I_{n-1}$ . We then extract the center  $p \times p$  patch from the flow map to create the flow feature vector for each contour.

### 3.2. Contour Merging

With the contour graph constructed, we observe that a single tree crown is often split into multiple contours. In other words, to predict accurate instance segmentation, one would need to merge these contours (Fig. 5). Hence, we formulate contour merging as an edge classification problem on a graph.

**Edge Classification.** For a set of contours  $\mathcal{C}$ , we represent whether pairs of contours belong to the same instance via a matrix  $\mathbf{Y} \in \{0, 1\}^{|\mathcal{C}| \times |\mathcal{C}|}$ , where  $Y_{ij} = 1$  indicates that contours  $C_i$  and  $C_j$  should be merged. Given a contour graph input, we aim to build a model that predicts this merge matrix. We formulate this task as an edge classification problem of the contour graph.

We have extracted various features associated with each node in the graph (Sec. 3.1). We concatenate these features into a node feature matrix  $\mathbf{H}^0 \in \mathbb{R}^{|\mathcal{C}| \times d}$ . Next, we represent the edges via an adjacency matrix  $\mathbf{A} \in [0, 1]^{|\mathcal{C}| \times |\mathcal{C}|}$ .

**Graph Convolution.** Using a graph convolution layer [29], we perform  $L$  rounds of message-passing to aggregate information from each node's neighborhood:

$$\mathbf{H}^{(l+1)} = f_{\text{GCN}}^{(l)}(\mathbf{H}^{(l)}) = \sigma(\tilde{\mathbf{A}}\mathbf{H}^{(l)}\mathbf{W}^{(l)}), \quad (1)$$

where  $\tilde{\mathbf{A}}$  denotes the degree normalized adjacency matrix,  $\sigma$  denotes an element-wise non-linearity, and  $\mathbf{W}^{(l)}$  denotes the trainable weights at the  $l^{\text{th}}$  layer.

To create the edge features for classification, following Mikolov et al. [41], we average the  $L^{\text{th}}$  node features  $\mathbf{H}^{(L)}$  for the pair of nodes in an edge  $(C_i, C_j)$ :

$$\mathbf{E}_{ij} = \frac{1}{2} \left( \mathbf{H}_i^{(L)} + \mathbf{H}_j^{(L)} \right). \quad (2)$$

Given the edge features  $\mathbf{E}_{ij}$ , we predict whether a pair of nodes should be merged using a multi-layer perceptron, i.e.,

$$\hat{Y}_{ij} = \text{MLP}(\mathbf{E}_{ij}) \quad \forall (C_i, C_j \in \mathcal{E}). \quad (3)$$

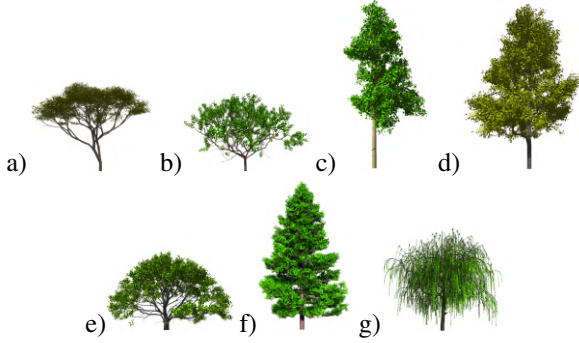


Figure 6. **Tree Models.** Examples of synthetic trees: a) acacia, b) apple, c) birch, d) maple, e) oak, f) pine, and g) willow).

To train this model, we minimize the binary cross-entropy loss between the predicted merge matrix  $\hat{\mathbf{Y}}$  and the ground truth merge matrix  $\mathbf{Y}$  over the edges:

$$\mathcal{L} = - \sum_{\{(i,j):A_{ij}=1\}} \mathbf{Y}_{ij} \log(\hat{\mathbf{Y}}_{ij}) + (1 - \mathbf{Y}_{ij}) \log(1 - \hat{\mathbf{Y}}_{ij}).$$

Finally, we merge the contours according to the predicted merge matrix  $\hat{\mathbf{Y}}$  and output an instance mask for each of the merged contours. In other words, for a tree  $t$ , the predicted mask is the union of all pixel locations within merged contours, i.e.,

$$\mathbf{M}_t = \bigcup_{j:\mathbf{Y}_{tj}=1} \mathcal{C}_j. \quad (4)$$

### 3.3. Collection of Datasets

We have prepared four datasets to develop and evaluate our method. We have made the synthetic and real-world dataset (with videos) publicly available. Please visit <https://github.com/adnan0819/Tree-Instance-Segmentation-using-Temporal-Structured-Images>.

**Synthetic Dataset.** A first dataset is created synthetically and enables us to control the density, species, and coloring of the trees. We use the algorithm for plant competition for resources [46] implemented as a developmental model [42] controlled by user-defined parameters [33, 57] and its realism has been validated by a prior large user study [48]. The trees have been placed by ecosystem simulation [7, 44]. Our implementation supports acacia, apple, willow, maple, birch, oak, and pine trees (see Fig. 6). We can generate terrains of arbitrary size full of trees.

**Real-world Datasets.** We have prepared three real-world forest datasets. Forest A (Martell Foret, West Lafayette, IN) was collected by our team with a UAV flying over a large forest and we made this dataset available publicly. To our knowledge, it is the largest and most comprehensive UAV-captured overhead image sequence dataset for forest research. Moreover, although not used here, each image includes georegistered coordinates of the camera location.

| Methods              | GT   | Pred. | Acc. $\uparrow$ |
|----------------------|------|-------|-----------------|
| YOLOv7               | 2172 | 3442  | 41.5            |
| YOLOv7 + Flow        | 2172 | 3134  | 55.7            |
| YOLOv7 + Flow (Med.) | 2172 | 3016  | 61.1            |
| YOLOv7 + Flow (Mean) | 2172 | 2838  | 69.3            |
| Ours                 | 2172 | 2373  | <b>90.7</b>     |

Table 1. **Preliminary experiments.** Count accuracy on synthetic forest verifying the effectiveness of capturing self-occlusion using optical flow. Med. and Mean correspond to thresholding the optical flow based on the median or mean magnitude.

The UAV was piloted by a trained forestry pilot and used two cameras: DJI P1 and DJI H20T. Three distinct forest regions were captured using 23 flights and covering approximately 368 hectares (3,680,000 m<sup>2</sup>) in total. The image sequences were captured at 80m, 100m, and 120m altitudes with the camera sensors pointing straight down (i.e., nadir). The DJI P1 camera had a field of view (FOV) of 63.5°, and DJI H20T had a FOV of 82.9°. All flight speeds were 5m/s. Each image has a resolution of 3,840 × 2,160 pixels and was captured at the rate of 60 images per second.

To evaluate the generalization capability (i.e., *out-of-distribution*), we collected Forest B (Kentucky Ridge State Forest, KY) and C (Olympic National Forest, WA) from Google Earth. We annotated twenty images at 832 × 732 pixels with tree crown information. Forest B and C are only used for validation. Notably, these forests have different characteristics. Forest A has a combination of natural and plantation vegetation and combines both deciduous and coniferous species like white oak, black cherry, red oak, etc. In contrast, Forest B contains more color variation without any plantation-type forest with species like sugar maple, tulip poplar, various oaks, hemlocks, etc. Differently, Forest C contains both seasonally changing and evergreen species. These differences in datasets B and C make them suitable for out-of-distribution evaluation.

## 4. Results

We evaluate our proposed approach using the Synthetic, Forest A, B, and C datasets. We report on the task of tree crown instance segmentation and counting. Note that instance segmentation systems trivially generalize to counting by predicting the number of detected instances. Beyond quantitative comparisons with baselines, we also conduct ablation studies verifying the efficacy of our proposed components and qualitative demonstrations.

**Experiment Setup.** For synthetic data experiments, we use 80% of the data to train and report evaluation metrics on the remaining 20%. In real-world experiments, we train on Forest A also using an 80-20 train and validation split. Using the same model trained on A, we report performance

| Methods            | Synthetic Forest |                             |                             | Forest A      |                             |                             | Forest B      |                             |                             | Forest C      |                             |                             |
|--------------------|------------------|-----------------------------|-----------------------------|---------------|-----------------------------|-----------------------------|---------------|-----------------------------|-----------------------------|---------------|-----------------------------|-----------------------------|
|                    | AP $\uparrow$    | AP <sub>50</sub> $\uparrow$ | AP <sub>70</sub> $\uparrow$ | AP $\uparrow$ | AP <sub>50</sub> $\uparrow$ | AP <sub>70</sub> $\uparrow$ | AP $\uparrow$ | AP <sub>50</sub> $\uparrow$ | AP <sub>70</sub> $\uparrow$ | AP $\uparrow$ | AP <sub>50</sub> $\uparrow$ | AP <sub>70</sub> $\uparrow$ |
| Mask-RCNN (ResNet) | 27.1             | 53.6                        | 50.1                        | 33.4          | 57.3                        | 54.1                        | 39.2          | 58.8                        | 56.1                        | 35.1          | 57.9                        | 58.4                        |
| Mask-RCNN (Swin-T) | 59.8             | 70.2                        | 68.3                        | 64.6          | 74.1                        | 70.5                        | 69.5          | <b>77.3</b>                 | <b>72.4</b>                 | 63.2          | 72.6                        | 70.3                        |
| TraDeS             | 61.1             | 55.2                        | 64.4                        | 58.1          | 71.3                        | 66.8                        | 63.7          | 73.9                        | 70.5                        | 59.6          | 70.4                        | 64.1                        |
| BoundaryFormer     | 56.3             | 65.1                        | 57.5                        | 60.9          | 72.9                        | 66.2                        | 64.1          | 73.4                        | 69.2                        | 58.9          | 71.2                        | 61.8                        |
| Mask2Former        | 62.4             | 65.2                        | 63.9                        | 59.7          | 64.1                        | 63.2                        | 64.1          | 67.5                        | 63.8                        | 61.9          | 63.1                        | 62.9                        |
| MS-RCNN            | 66.2             | 68.4                        | 65.8                        | 64.8          | 71.3                        | 68.5                        | 66.4          | 70.2                        | 69.3                        | 64.2          | 69.3                        | 65.8                        |
| OCISIS             | 59.1             | 63.5                        | 61.8                        | 55.7          | 58.6                        | 58.1                        | 60.8          | 68.2                        | 66.8                        | 60.2          | 62.8                        | 60.9                        |
| SLIC (Superpixel)  | 23.7             | 27.3                        | 24.2                        | 20.7          | 24.8                        | 23.6                        | 22.5          | 27.3                        | 23.7                        | 20.1          | 27.5                        | 24.6                        |
| Aerial Laser       | 71.1             | 72.9                        | 70.8                        | 70.4          | 75.2                        | 71.1                        | 65.1          | 68.2                        | 66.8                        | 65.3          | 68.7                        | 64.8                        |
| Ours               | <b>74.6</b>      | <b>73.1</b>                 | <b>69.5</b>                 | <b>74.5</b>   | <b>81.6</b>                 | <b>72.8</b>                 | <b>69.8</b>   | 76.2                        | 71.5                        | <b>70.1</b>   | <b>75.4</b>                 | <b>72.5</b>                 |

Table 2. **Segmentation.** Comparisons of segmentation performance to prior instance segmentation baselines.

| Methods            | Synthetic Forest |       |                 | Forest A |       |                 | Forest B |       |                 | Forest C |                  |                 |
|--------------------|------------------|-------|-----------------|----------|-------|-----------------|----------|-------|-----------------|----------|------------------|-----------------|
|                    | GT               | Pred. | Acc. $\uparrow$ | GT       | Pred. | Acc. $\uparrow$ | GT       | Pred. | Acc. $\uparrow$ | GT       | Pred. $\uparrow$ | Acc. $\uparrow$ |
| Mask-RCNN (ResNet) | 5157             | 3291  | 63.8            | 2314     | 1691  | 73.1            | 2041     | 1281  | 62.8            | 2172     | 1403             | 64.6            |
| Mask-RCNN (Swin-T) | 5157             | 4275  | 82.9            | 2314     | 1816  | 78.5            | 2041     | 1655  | 81.1            | 2172     | 1791             | 82.5            |
| TraDeS             | 5157             | 4131  | 80.1            | 2314     | 1987  | 85.9            | 2041     | 1596  | 78.2            | 2172     | 1611             | 74.2            |
| BoundaryFormer     | 5157             | 3981  | 77.2            | 2314     | 1950  | 84.3            | 2041     | 1549  | 75.9            | 2172     | 1713             | 78.9            |
| Mask2Former        | 5157             | 3290  | 63.8            | 2314     | 1708  | 73.8            | 2041     | 1601  | 78.4            | 2172     | 1681             | 77.4            |
| MS-RCNN            | 5157             | 3527  | 68.4            | 2314     | 1886  | 81.6            | 2041     | 1645  | 80.6            | 2172     | 1752             | 80.7            |
| OCISIS             | 5157             | 4373  | 84.8            | 2314     | 1969  | 85.1            | 2041     | 1643  | 80.5            | 2172     | 1822             | 83.9            |
| SLIC (Superpixel)  | 5157             | 8142  | 42.1            | 2314     | 3566  | 45.9            | 2041     | 3298  | 38.4            | 2172     | 3521             | 37.9            |
| Aerial Laser       | 5157             | 4399  | 85.3            | 2314     | 1960  | 84.7            | 2041     | 1639  | 80.3            | 2172     | 1887             | 86.9            |
| Ours               | 5157             | 5636  | <b>90.7</b>     | 2314     | 2510  | <b>91.5</b>     | 2041     | 1771  | <b>86.8</b>     | 2172     | 2384             | <b>90.2</b>     |

Table 3. **Counts.** Comparisons of tree crown count performance to prior instance segmentation baselines.

on B and C to evaluate out-of-distribution generalization.

**Evaluation Metrics.** We follow the evaluation protocol from the prior work [26]. For *tree instance segmentation*, we report the mask Average Precision (AP) averaged over different intersection-over-union (IoU) thresholds and at IoU thresholds of 0.5 and 0.7, denoted as AP<sub>50</sub> and AP<sub>70</sub> respectively. For *tree counting*, we report the raw number of counts and the count accuracy, which is expressed as a percentage denoted by Acc.

**Baselines.** We compare our approach to the following instance segmentation methods: Mask-RCNN based on ResNet [26], Mask-RCNN based on the recent transformer backbone (Swin-T) [40], Track-to-Detect-and-Segment (TraDeS) [59], BoundaryFormer [32] (a recent mask-supervised polygonal boundary approach to instance segmentation using transformers), Mask2Former [14], Mask Scoring RCNN [27], Object Counting and Instance Segmentation with Image-level Supervision [15], SLIC (superpixels) [3], and Aerial Scanning Using Laser and Deep Model (a model specific to trees) [54].

#### 4.1. Synthetic Data Experiments

**Preliminary Experiments.** We used our synthetic data to aid the development of our method. In particular, we used it to understand the impact of determining the change in self-

occlusion, and we experimented with several settings using our synthetic dataset.

First, we varied the density of the tree crowns in the dataset until YOLOv7 [58] had difficulty in determining the tree crown count (see Tab. 1).

Second, we used the graphics rendering engine to determine the ground truth optical flow as a virtual UAV flies over the same dense synthetic forest. We extended YOLOv7 [58] to include this flow data as an indicator of changing self-occlusion patterns. Specifically, we incorporate flow as additional channels to the input.

We found that different optical flow thresholds led to improved count accuracies. Hence, Tab. 1 reports the counting performance using flow when not thresholded, thresholded by the median and the mean. The incorporation of the optical flow improves the performance by 33%. However, the counting performance of YOLOv7 remains unsatisfactory.

To further improve performance, we developed our method to include pixel content features, shape features and changed the underlying model to the graph-based approach described in Sec. 3 – this produced the highest count accuracy of 90.7%, as seen in Tab. 1. In contrast, the next best model, YOLOv7+Flow(mean), achieves a count accuracy of only 69.3%.

**Quantitative Results.** We compared the count accuracy

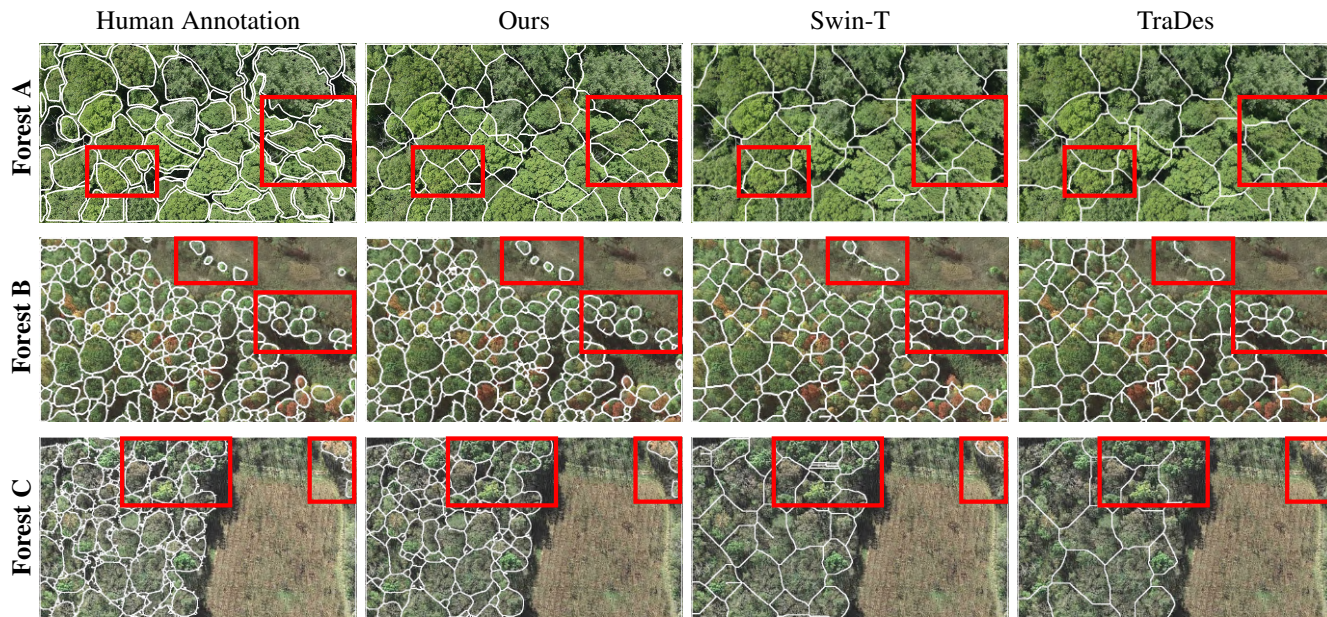


Figure 7. **Qualitative comparisons.** Comparisons of our approach vs. Swin-T and TraDes. We observe that our approach produces instance segmentation masks that more closely match the human annotation. Comparisons to more baselines are illustrated in the Appendix.

of our method using the synthetic dataset to our baselines in Tab. 3 (first three columns). We observe that MaskRCNN (Swin-T) is the best-performing baseline with a count accuracy of 82.9%, followed by TraDeS with 80.1%, BoundaryFormer with 77.2%, and Mask-RCNN (ResNet) with 63.8%. Our method outperforms all baselines in count accuracy with 90.7%. We also compared the AP performance of using the synthetic dataset to our baselines in Tab. 2. Our method again outperforms all the baselines.

## 4.2. Real-World Data Experiments

We also conducted experiments on our multiple real-world forest datasets. Tab. 2 contains the AP comparisons to the baselines. Our approach outperforms all methods in all forests except for AP<sub>50</sub> and AP<sub>70</sub> in Forest B. Overall, our method achieves an AP of 74.5, 69.8 and 70.1 on Forest A, B, and C respectively.

Tab. 3 reports the counting accuracy comparison to our baselines. The accuracy metric is defined as  $1 - MAE$  (normalized) as a percentage where MAE is the mean absolute error. Our method exceeds the performance of all prior methods across all forests. On Forest A, our method achieves a count accuracy of 91.5% compared to the next best baseline (TraDeS) of 85.9%. On Forest B, our method achieves 86.8 compared to the next best baseline (Swin-T) of 81.1%. On Forest C, our method achieves 90.2% compared to the next best baseline (Aerial Laser Scanning Model) of 86.9%.

| Ablated feature      | AP   | AP <sub>50</sub> | AP <sub>75</sub> | Acc. |
|----------------------|------|------------------|------------------|------|
| All Features         | 74.5 | 81.6             | 72.8             | 91.5 |
| All — aspect ratio   | 67.6 | 72.9             | 68.4             | 82.4 |
| All — solidity       | 69.2 | 74.1             | 66.2             | 80.1 |
| All — self-occlusion | 62.9 | 77.3             | 71.1             | 78.2 |
| All — patch          | 51.3 | 56.8             | 63.2             | 71.8 |
| All — deviation      | 59.8 | 62.2             | 64.1             | 70.9 |
| All — area           | 55.1 | 58.8             | 61.6             | 68.3 |
| All — neighbor sim.  | 57.4 | 61.7             | 59.6             | 63.5 |

Table 4. **Feature ablations.** Performance of our approach on Forest A as we ablate each feature (sorted by count accuracy). Features with lower accuracy have a higher impact on the prediction.

## 4.3. Qualitative Study

Beyond the quantitative results, we also qualitatively study the method’s behavior. Fig. 7 shows the results of our method and baselines applied to the three real-world datasets. We generally observe that the predicted instance segmentation closely matches the ground truth annotations. On the other hand, we observe that a typical failure case for TraDes and Swin-T is that they often group multiple tree crowns into one; See Fig. 7 for regions boxed in red. —Similar illustrations of qualitative comparisons with more baselines are given in the Appendix.

Next, we study the merging behavior (Sec. 3.2). Fig. 8 shows the contour map before and after merging. We observed that our approach merged small contours into larger ones to match the tree crown boundaries more closely.

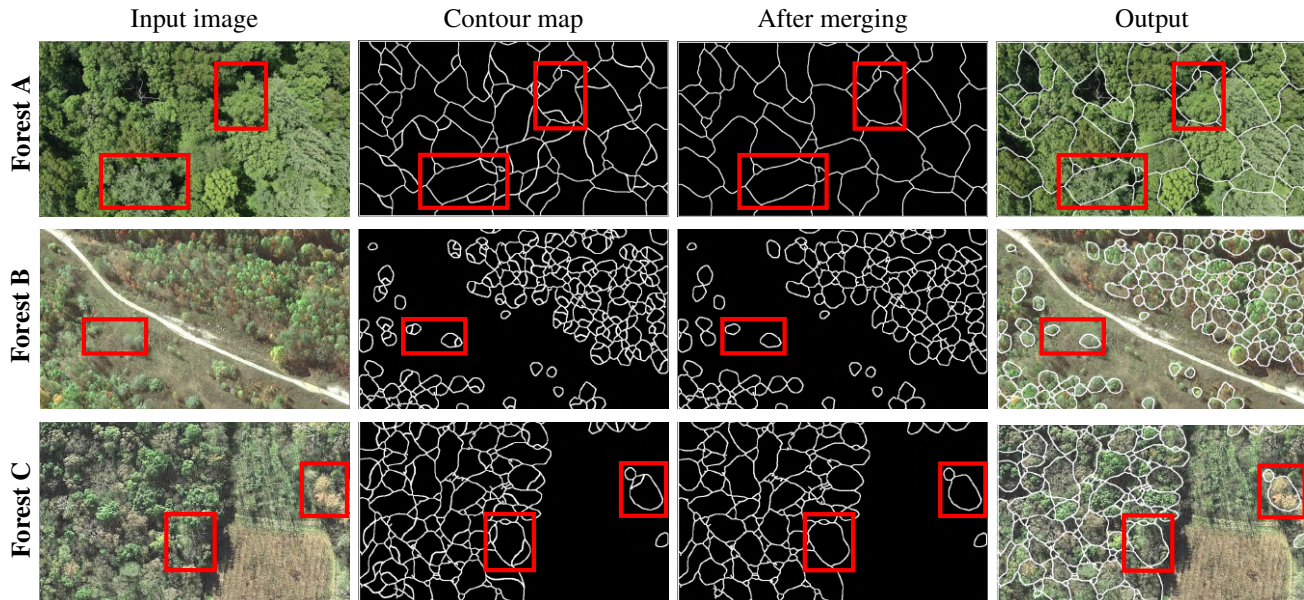


Figure 8. **Visual illustration of data flow.** We visualize the contour map before and after contour merging.

| Frame $\Delta$ | 1    | 2    | 3    | 4    | 5    | no flow |
|----------------|------|------|------|------|------|---------|
| AP $\uparrow$  | 74.5 | 72.8 | 68.1 | 65.4 | 64.8 | 62.9    |
| Acc $\uparrow$ | 91.5 | 86.1 | 83.5 | 82.9 | 81.8 | 78.2    |

Table 5. **Frame-rate ablations.** The performance drop as optical flow is computed at different image intervals.

#### 4.4. Ablation Study

To quantify the importance of features of our contour graph, we perform an ablation study using Forest A. In Tab. 4, we sort the ablations by descending in count accuracy – i.e., the feature in the top rows are the features that least impact model performance. We observe that ‘neighbor similarity’ impacts count accuracy the most, and the patch feature has the most impact on AP.

The other shape and self-occlusion features contribute significantly to the segmentation. The ablation shows a 13.3% increase in count accuracy when the self-occlusion features are used (driven by optical flow). Not using this feature in dense scale forests, similar to the samples of all real forest trees reported in Tab. 3, would lead to an over or underestimation of approximately 868 trees.

Moreover, to test the capability of the optical flow algorithm used [17], we skipped progressively more frames and recomputed the flow in each case (ablating the self-occlusion feature). The different flow values were then used and evaluated in terms of AP (see Tab. 5) and tree count accuracy. We can see in Tab. 5 the largest shift happens from skipping one frame to two frames, as the granularity of flow (using motion) halves. The subsequent drops indicate lower

frame rate adversely affects the performance. Notably, even at lower frame rates, the performance metrics were higher than having no optical flow. In summary, we observe that our method generally functions well even for images captured between 30-60 fps.

## 5. Limitations & Conclusion

Our method is limited by the quality of the input image. For example, trees include very high-frequency details, e.g., small branchlets, that are below the resolution of the UAV cameras, which our approach cannot capture. Other image degradation factors, such as blurring caused by the wind and the UAV motion, also lead to similar challenges. Also, some shadows can be identified as trees. An obvious limitation is that our approach does not work for a single image and requires an image sequence.

In summary, we introduced a novel approach to instance segmentation and counting tree crowns in forests. Our approach uses image sequences from increasingly available UAVs. Our key contribution is leveraging the partial occlusion between successive images, shape features of the contours, and encoding these features into a contour graph that is updated between successive images. The result is a method that produces instance segmentation masks that separate the individual tree crowns. Finally, we provide synthetic and real-world datasets that can facilitate future research in tree crown instance segmentation and counting.

**Acknowledgements:** This research is at least partially supported by Purdue Digital Forestry Center (esp. Dr. Songlin Fei) and NSF grants 1835739 and 2106717.



## References

- [1] Real time design and animation of fractal plants and trees. *ACM SIGGRAPH Comput. Graph.*, 1986. 3
- [2] Review on convolutional neural networks (CNN) in vegetation remote sensing. *ISPRS P&RS*, 2021. 3
- [3] Radhakrishna Achanta, Appu Shaji, Kevin Smith, Aurelien Lucchi, Pascal Fua, and Sabine Süsstrunk. SLIC superpixels compared to state-of-the-art superpixel methods. *TPAMI*, 2012. 6
- [4] M. Aono and T.L. Kunii. Botanical tree image generation. *IEEE Computer Graphics and Applications*, 1984. 3
- [5] Hasan Asy'ari Arief, Geir-Harald Strand, Håvard Tveite, and Ulf Geir Indahl. Land cover segmentation of airborne lidar data using stochastic atrous network. *Remote Sensing*, 2018. 1
- [6] James Arvo and David Kirk. Modeling plants with environment-sensitive automata. In *Proc. of Ausgraph*, 1988. 3
- [7] Bedrich Benes, Nathan Andryscio, and Ondřej Št'ava. Interactive modeling of virtual ecosystems. In *Proc. Eurographics Conference on Natural Phenomena*, 2009. 5
- [8] Bedrich Benes and Erik Uriel Millán. Virtual climbing plants competing for space. In *Computer Animation*. IEEE Computer Society, 2002. 3
- [9] Alexey Bochkovskiy, Chien-Yao Wang, and Hong-Yuan Mark Liao. YOLOv4: Optimal speed and accuracy of object detection. *arXiv preprint arXiv:2004.10934*, 2020. 2
- [10] Daniel Bolya, Chong Zhou, Fanyi Xiao, and Yong Jae Lee. Yolact: Real-time instance segmentation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019. 3
- [11] Liang-Chieh Chen, Alexander Hermans, George Papandreou, Florian Schroff, Peng Wang, and Hartwig Adam. MaskLab: Instance segmentation by refining object detection with semantic and direction features. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2018. 3
- [12] Steven W. Chen, Guilherme V. Nardari, Elijah S. Lee, Chao Qu, Xu Liu, Roseli Ap. Francelin Romero, and Vijay R. Kumar. Sloam: Semantic lidar odometry and mapping for forest inventory. *IEEE RA-L*, 2020. 3
- [13] Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022. 3
- [14] Bowen Cheng, Ishan Misra, Alexander G. Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In *CVPR*, 2022. 6
- [15] Hisham Cholakkal, Guolei Sun, Fahad Shahbaz Khan, and Ling Shao. Object counting and instance segmentation with image-level supervision. In *CVPR*, 2019. 6
- [16] Phillippe De Reffye, Claude Edelin, Jean Françon, Marc Jaeger, and Claude Puech. Plant models faithful to botanical structure and development. *ACM Siggraph Comp. Graph.*, 1988. 3
- [17] Gunnar Farneback. Two-frame motion estimation based on polynomial expansion. In Josef Bigun and Tomas Gustavsson, editors, *Image Analysis*, 2003. 4, 8
- [18] Adnan Firoze, Bedrich Benes, and Daniel Aliaga. Urban tree generator: spatio-temporal and generative deep learning for urban tree localization and modeling. *The Visual Computer*, 2022. 2
- [19] Ross Girshick. Fast R-CNN. In *Int. Conf. Comput. Vis.*, 2015. 2
- [20] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2014. 3
- [21] Ned Greene. Voxel space automata: Modeling with stochastic growth processes in voxel space. *ACM SIGGRAPH Comp. Graph.*, 1989. 3
- [22] Zicheng Guo and Richard W. Hall. Parallel thinning with two-subiteration algorithms. *ACM Commun.*, 1989. 4
- [23] Torsten Hädrich, Daniel T. Banuti, Wojtek Pałubicki, Sören Pirk, and Dominik L. Michels. Fire in paradise: Mesoscale simulation of wildfires. *ACM Trans. Graph.*, 2021. 3
- [24] Torsten Hädrich, Bedrich Benes, Oliver Deussen, and Sören Pirk. Interactive modeling and authoring of climbing plants. 2017. 3
- [25] Bharath Hariharan, Pablo Arbeláez, Ross Girshick, and Jitendra Malik. Simultaneous detection and segmentation. In *Eur. Conf. Comput. Vis.*, 2014. 3
- [26] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask R-CNN. In *Int. Conf. Comput. Vis.*, 2017. 1, 3, 6
- [27] Zhaojin Huang, Lichao Huang, Yongchao Gong, Chang Huang, and Xinggang Wang. Mask scoring R-CNN. In *CVPR*, 2019. 6
- [28] Lei Ke, Yu-Wing Tai, and Chi-Keung Tang. Deep occlusion-aware instance segmentation with overlapping bilayers. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021. 1, 3
- [29] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *Int. Conf. Learn. Represent.*, 2017. 4
- [30] Alexander Kirillov, Yuxin Wu, Kaiming He, and Ross Girshick. PointRend: Image segmentation as rendering. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020. 3
- [31] Hei Law and Jia Deng. Cornernet: Detecting objects as paired keypoints. In *Eur. Conf. Comput. Vis.*, 2018. 2
- [32] Justin Lazarow, Weijian Xu, and Zhuowen Tu. Instance segmentation with mask-supervised polygonal boundary transformers. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022. 3, 6
- [33] Bosheng Li, Jacek Kałużny, Jonathan Klein, Dominik L. Michels, Wojtek Pałubicki, Bedrich Benes, and Sören Pirk. Learning to reconstruct botanical trees from single images. *ACM Trans. Graph.*, 2021. 2, 5
- [34] Chuyi Li, Lulu Li, Hongliang Jiang, Kaiheng Weng, Yifei Geng, Liang Li, Zaidan Ke, Qingyuan Li, Meng Cheng, Weiqiang Nie, Yiduo Li, Bo Zhang, Yufei Liang, Linyuan Zhou, Xiaoming Xu, Xiangxiang Chu, Xiaoming Wei, and Xiaolin Wei. Yolov6: A single-stage object detection framework for industrial applications, 2022. 2
- [35] Dongze Lian, Jing Li, Jia Zheng, Weixin Luo, and Shenghua Gao. Density map regression guided detection network for rgb-d crowd counting and localization. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019. 1, 2
- [36] Lingbo Liu, Jiaqi Chen, Hefeng Wu, Guanbin Li, Chenglong Li, and Liang Lin. Cross-modal collaborative representation learning and a large-scale rgbt benchmark for crowd count-

- ing. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021. [1](#), [2](#)
- [37] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. SSD: Single shot multibox detector. In *Eur. Conf. Comput. Vis.*, 2016. [2](#)
- [38] Xinni Liu, Fengrong Han, Kamarul Hawari Ghazali, Izzeldin Ibrahim Mohamed, and Yue Zhao. A review of convolutional neural networks in remote sensing image. In *IC-SCA*, 2019. [3](#)
- [39] Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie, Yixuan Wei, Jia Ning, Yue Cao, Zheng Zhang, Li Dong, et al. Swin transformer v2: Scaling up capacity and resolution. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022. [3](#)
- [40] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Int. Conf. Comput. Vis.*, 2021. [3](#), [6](#)
- [41] Tomas Mikolov, Kai Chen, Gregory S. Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. In *Int. Conf. Learn. Represent.*, 2013. [4](#)
- [42] Radomír Měch and Przemyslaw Prusinkiewicz. Visual models of plants interacting with their environment. In *ACM SIG-GRAPH Comp. Graph.*, 1996. [3](#), [5](#)
- [43] Alejandro Newell and Jia Deng. Pixels to graphs by associative embedding. In *Adv. Neural Inform. Process. Syst.*, 2017. [1](#)
- [44] Till Niese, Sören Pirk, Matthias Albrecht, Bedrich Benes, and Oliver Deussen. Procedural urban forestry. *ACM Trans. Graph.*, 2022. [3](#), [5](#)
- [45] Masanori Onishi and Takeshi Ise. Explainable identification and mapping of trees using uav rgb image and deep learning. *Nature: Scientific Reports*, 2021. [3](#)
- [46] Wojciech Palubicki, Kipp Horel, Steven Longay, Adam Runions, Brendan Lane, Radomír Měch, and Przemyslaw Prusinkiewicz. Self-organizing tree models for image synthesis. *ACM Trans. Graph.*, 2009. [3](#), [5](#)
- [47] Sören Pirk, Till Niese, Torsten Hädrich, Bedrich Benes, and Oliver Deussen. Windy trees: Computing stress response for developmental tree models. *ACM Trans. Graph.*, 2014. [3](#)
- [48] Tomas Polasek, David Hrusa, Bedrich Benes, and Martin Cadik. Ictree: Automatic perceptual metrics for tree models. *ACM Trans. Graph.*, 2021. [3](#), [5](#)
- [49] Przemyslaw Prusinkiewicz and Aristid Lindenmayer. *The Algorithmic Beauty of Plants*. 1990. [3](#)
- [50] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2016. [2](#)
- [51] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *Adv. Neural Inform. Process. Syst.*, 2015. [2](#)
- [52] Alvy Ray Smith. Plants, fractals, and formal languages. In *ACM SIGGRAPH Comp. Graph.*, 1984. [3](#)
- [53] Z. Su, Wenzhe Liu, Zitong Yu, Dewen Hu, Qing Liao, Qi Tian, Matti Pietikäinen, and Li Liu. Pixel difference networks for efficient edge detection. In *Int. Conf. Comput. Vis.*, 2021. [4](#)
- [54] Chenxin Sun, Chengwei Huang, Huaqing 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*, 2022. [6](#)
- [55] Mingxing Tan, Ruoming Pang, and Quoc V Le. Efficient-Det: Scalable and efficient object detection. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020. [2](#)
- [56] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. Fcos: Fully convolutional one-stage object detection. In *Int. Conf. Comput. Vis.*, 2019. [2](#)
- [57] Ondrej Štáva, Sören Pirk, Julian Kratt, Baoquan Chen, Radomír Měch, Oliver Deussen, and Bedrich Benes. Inverse procedural modelling of trees. *Comput. Graph. Forum*, 2014. [2](#), [5](#)
- [58] Chien-Yao Wang, Alexey Bochkovskiy, and Hong-Yuan Mark Liao. Yolov7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors. *arXiv preprint arXiv:2207.02696*, 2022. [2](#), [6](#)
- [59] Jialian Wu, Jiale Cao, Liangchen Song, Yu Wang, Ming Yang, and Junsong Yuan. Track to detect and segment: An online multi-object tracker. In *Int. Conf. Comput. Vis.*, 2021. [3](#), [6](#)
- [60] Jianwei Yang, Chunyuan Li, Pengchuan Zhang, Xiyang Dai, Bin Xiao, Lu Yuan, and Jianfeng Gao. Focal attention for long-range interactions in vision transformers. In *Adv. Neural Inform. Process. Syst.*, 2021. [3](#)
- [61] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2018. [4](#)
- [62] Xizhou Zhu, Han Hu, Stephen Lin, and Jifeng Dai. Deformable ConvNetsv2: More deformable, better results. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019. [3](#)