



Contents lists available at ScienceDirect

# International Journal of Applied Earth Observation and Geoinformation

journal homepage: [www.elsevier.com/locate/jag](http://www.elsevier.com/locate/jag)

## Individual tree detection in large-scale urban environments using high-resolution multispectral imagery

Jonathan Ventura<sup>a,\*</sup>, Camille Pawlak<sup>b</sup>, Milo Honsberger<sup>c</sup>, Cameron Gonsalves<sup>c</sup>, Julian Rice<sup>a</sup>, Natalie L.R. Love<sup>b</sup>, Skyler Han<sup>a</sup>, Viet Nguyen<sup>a</sup>, Keilana Sugano<sup>c,d</sup>, Jacqueline Doremus<sup>e</sup>, G. Andrew Fricker<sup>c</sup>, Jenn Yost<sup>b</sup>, Matt Ritter<sup>b</sup>

<sup>a</sup> Department of Computer Science & Software Engineering, California Polytechnic State University, 1 Grand Ave, San Luis Obispo, 93107-0354, CA, USA

<sup>b</sup> Department of Biological Sciences, California Polytechnic State University, 1 Grand Ave, San Luis Obispo, 93107-0401, CA, USA

<sup>c</sup> Department of Social Sciences, California Polytechnic State University, 1 Grand Ave, San Luis Obispo, 93107-0329, CA, USA

<sup>d</sup> Department of Political Science, California Polytechnic State University, 1 Grand Ave, San Luis Obispo, 93107-0329, CA, USA

<sup>e</sup> Orfalea College of Business, California Polytechnic State University, 1 Grand Ave, San Luis Obispo, 93107, CA, USA

### ARTICLE INFO

Dataset link: <https://github.com/jonathanventura/urban-tree-detection>, <https://github.com/jonathanventura/urban-tree-detection-data>, <https://ufe.calpoly.edu/>

#### Keywords:

Tree detection  
Urban forests  
Multispectral imagery  
Aerial imagery  
Computer vision  
Object detection  
Deep learning  
Convolutional neural network

### ABSTRACT

Systematic maps of urban forests are useful for regional planners and ecologists to understand the spatial distribution of trees in cities. However, manually-created urban forest inventories are expensive and time-consuming to create and typically do not provide coverage of private land. Toward the goal of automating urban forest inventory through machine learning techniques, we performed a comparative study of methods for automatically detecting and localizing trees in multispectral aerial imagery of urban environments, and introduce a novel method based on convolutional neural network regression. Our evaluation is supported by a new dataset of over 1,500 images and almost 100,000 tree annotations, covering eight cities, six climate zones, and three image capture years. Our method outperforms previous methods, achieving 73.6% precision and 73.3% recall when trained and tested in Southern California, and 76.5% precision 72.0% recall when trained and tested across the entire state. To demonstrate the scalability of the technique, we produced the first map of trees across the entire urban forest of California. The map we produced provides important data for the planning and management of California's urban forest, and establishes a proven methodology for potentially producing similar maps nationally and globally in the future.

### 1. Introduction

Urban forests provide extensive benefits to residents of cities such as contributing to resident energy-savings, reducing impervious runoff and water quality, controlling microclimate, and sequestering carbon (McPherson and Simpson, 2002; McPherson et al., 2016; Livesley et al., 2016). City, state, or county governments manage the public section in order to estimate and maximize those benefits and make sure they are equitably distributed throughout a city. Governments typically track their urban forests through manual tree inventories performed by professional arborists (Nielsen et al., 2014). However, such inventories are typically limited to information about public street trees, which make up only a portion of a city's urban forest. The extent of the benefits provided to residents by their urban trees depends on the number of trees in that city, both public and private. When making policy decisions about where the next tree planting is needed most, cities need to be able to account for both the publicly and

privately managed urban forest in order to target plantings. Automatic tree detection implemented using machine learning and aerial imagery provides a method for cities to efficiently and cost-effectively track their urban forests.

Tree detection methods can be divided into two categories: methods that operate on LiDAR-derived products such as a canopy height model (CHM) (Chen and Zakhor, 2009; Silva et al., 2016; Liu et al., 2015; Wu et al., 2016; Zörner et al., 2018; Roussel et al., 2020; Xu et al., 2021; Münzinger et al., 2022) or a 3D point cloud (Ayrey and Hayes, 2018; Chen et al., 2021), and methods that operate solely on optical imagery. The geometric information obtained from LiDAR is highly useful for tree detection. However, LiDAR is expensive to collect, and there is currently no source for high-resolution LiDAR data for much of the Earth, unlike high-resolution multi-spectral imagery, for which multiple sources exist.

\* Corresponding author.

E-mail address: [jventu09@calpoly.edu](mailto:jventu09@calpoly.edu) (J. Ventura).

<https://doi.org/10.1016/j.jag.2024.103848>

Received 11 August 2023; Received in revised form 26 March 2024; Accepted 17 April 2024

Available online 2 May 2024

1569-8432/© 2024 Published by Elsevier B.V. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

The advantages of using optical, multispectral imagery are that it is collected everywhere in the United States as part of the U.S. National Agricultural Imagery Program (NAIP) and it is cheaper to collect than LiDAR. LiDAR has the benefit of precisely measuring vertical dimensions, but coverage is spotty with unreliable repeat times. The primary motivation for using NAIP imagery is its wall-to-wall coverage and repeated collections which would allow for tracking changes over time. Indeed, a blended approach using optical and LiDAR data would be preferable, but is not repeatable at scale.

Early methods for tree detection in remote sensing imagery perform per-pixel classification to identify tree canopies, tree characteristics, and tree species (Xiao et al., 2004; Yang et al., 2009; Jensen et al., 2012; Alonzo et al., 2013; Shang and Chisholm, 2014; Bosch, 2020; Brandt et al., 2020). The limitation of most of these approaches is that they output per-pixel maps rather than detecting and localizing individual trees. However, some methods have been proposed to extract individual tree locations from the per-pixel maps as a post-processing step, such as template matching (Yang et al., 2009) and connected components (Brandt et al., 2020).

Many recent methods on tree detection adopt an object detection approach (Santos et al., 2019; Zamboni et al., 2021; Weinstein et al., 2019; Zhang et al., 2022; Das et al., 2022; Beloiu et al., 2023) or an instance segmentation approach (Martins et al., 2021; Freudenberg et al., 2022; Yang et al., 2022; Sun et al., 2022; Ball et al., 2023). Object detection requires a bounding box annotation for each tree in the training dataset, while instance segmentation requires an accurate delineation of the tree crown. Bounding box and crown delineation annotations have size information which provides a useful supervisory signal for machine learning methods. However, fully annotating a bounding box or complete crown delineation for each tree is time-consuming, and it is often difficult for human annotators to visually determine the correct delineation of overlapping tree crowns.

In this study we focus on the alternative approach of annotating and predicting tree location points. This approach makes the annotation process easier and more scalable. However, object detection and instance segmentation methods are unable to learn from point annotations alone because they lack size information. Instead, we can consider methods from the substantial literature on object counting (Lempitsky and Zisserman, 2010), where the goal is to learn to count the number of objects in an image from only point annotations. Some previous work (Osco et al., 2020; Chen and Shang, 2022) has been successful in adapting these techniques for tree detection. In this study we expand upon the method of Osco et al. (2020) and introduce several modifications to improve the results.

Since our method only outputs point locations for the trees, our results cannot be used to directly estimate tree coverage area or tree crown size. However, the tree count estimates produced by our method would be generally useful to city managers, urban foresters, and ecologists. For example, when developing management plans, we believe city managers value tree counts over canopy estimates for estimating maintenance hours and staff requirements (Dwyer et al., 1992). Furthermore, ecologists could use our tree count estimates to track patterns in the urban forest over space and time (Love et al., 2022).

The primary objectives of our research are to produce a methodology for tree detection in urban areas in California and create a map of tree locations for every urban tree in California. In this work, we explore the potential for automatic detection and localization of trees in the California urban forest from aerial imagery. Using machine learning methods to detect trees in imagery, we can automatically produce a map of the trees in an urban environment, covering both public and private land and providing geographic coordinates for each tree.

To the best of our knowledge, no previous study has evaluated the potential for machine learning-based methods to automate urban tree detection across the entire state of California. Previous studies have either evaluated urban tree detection in individual cities in California (Wegner et al., 2016; Branson et al., 2018) or have addressed

related but different tasks, such as tree detection in the natural forests of California (Weinstein et al., 2019; Chen and Shang, 2022) or automated tree species identification from aerial imagery of California cities (Beery et al., 2022). Our dataset covers several cities across California, a state with considerable topographic and climatic diversity, and also includes imagery from multiple years to evaluate the potential for longitudinal analysis. Specifically, our hand-annotated dataset spans eight cities, six climate zones, and three image capture years, and including almost 100,000 tree annotations in total. The high value, broad impact, and diversity of California's urban forest make it such that studies done in this region are broadly applicable to urban forests globally.

## 2. Study area and dataset

### 2.1. Study area

California is the third largest state in the U.S. and the most populous with a population of 39,538,223 according to the 2020 Decennial Census (U.S. Census Bureau, 2020). The California urban forest provides services and benefits to California's residents previously estimated to be worth \$8.3 billion annually (McPherson et al., 2017). California's urban forests are composed of species with wide-ranging climatic requirements (Love et al., 2022). California's cities can host these species because the state has substantial topographic and climatic diversity, with six distinct California climate zones covering deserts, coastal and alpine environments (Love et al., 2022; Clarke et al., 2013). The California urban forest has high diversity compared to other urban forests, with a higher diversity ranking than the United States at a national scale, and a higher ranking relative to the most diverse urban forests within the United States, the Western region (Love et al., 2022; Ma et al., 2020). California's urban forests are so diverse that they are comparable not just to urban forests, but to the most diverse forests in the world, tropical forests (Love et al., 2022).

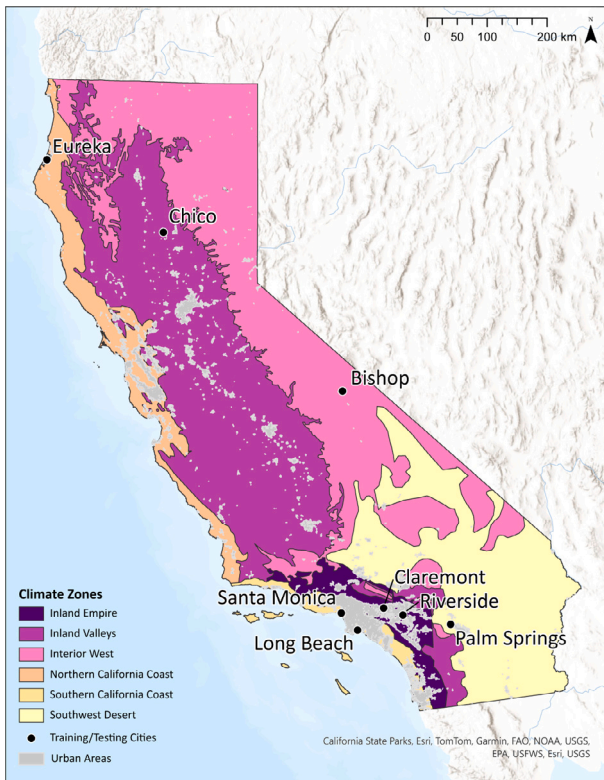
### 2.2. Dataset

We prepared our training and testing datasets using imagery from U.S. National Agricultural Program (NAIP), which contains aerial imagery acquired during the agricultural growing seasons in the United States. NAIP multispectral imagery is acquired every two years and covers the entire contiguous United States, typically at 1 m resolution or 60 cm resolution. We chose NAIP imagery because its coverage, temporal and spatial resolution, and public availability make it an ideal resource to support large-scale study of the urban forest in the United States. In our study we only used 60 cm imagery collected in 2016 or later.

Table 1 summarizes our dataset and Fig. 1 illustrates the locations of the included cities and climate zones. In total, our dataset contains 1651 images and 96,425 annotated trees, and covers eight cities and all six climate zones in California (McPherson, 2010).

The Southern California 2020 portion of the dataset covers five cities in Southern California and contains roughly 90–100 square image tiles captured in 2020 from each city: Claremont, Long Beach, Palm Springs, Riverside, and Santa Monica. Each image has a size of  $256 \times 256$  pixels, a resolution of 60 cm, and includes red, green, blue, and near infrared channels. We collected and annotated NAIP imagery from 2020 for these sites. We withheld a random 10% split of the Southern California 2020 images for testing and used the remaining 90% for training.

To further test the extrapolation ability of the network, we prepared two additional subsets of imagery. We annotated the same sites in Southern California in the Southern California 2020 portion of the dataset, but using imagery from 2016 and 2018, to test the ability of our method to process imagery of a similar location but captured at different times. We also selected and annotated images from three



**Fig. 1.** Locations of cities from which we collected and annotated images to form our dataset. Each climate zone in California is represented by at least one city in the dataset. The climate zones represent the following proportions of the total urban reserve area across the state: Inland Empire (24.03%), Inland Valleys (26.36%), Interior West (2.25%), Northern California Coast (16.46%), Southern California Coast (24.65%), and Southwest Desert (6.25%). (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

cities in Northern California: Bishop, Chico, and Eureka. These cities are in different climate zones from the Southern California sites, and in some places have a higher density of trees, since the cities include more densely forested rural areas. Thus, these images serve to test the ability of our method to extrapolate to other regions different from the training sites.

Our complete dataset combines the Southern and Northern California subsets and data from all three years (2016, 2018, and 2020). The dataset covers all five of California’s climate zones, with an emphasis on the Southern California Coast and Inland Empire, which contain 48.68% of the state’s total urban area. After developing our initial model using the Southern California 2020 subset and validating its extrapolation ability with the remaining images, we re-trained the model using the complete dataset to obtain our final tree detection model.

The various dataset subsets used in the paper are summarized in Table 2, which lists the number of images and tree annotations in the train and test splits for each subset.

We used tree inventories acquired by the cities as a starting point to annotate by hand the locations of all trees in the images, including trees in both public and private spaces. Since the city inventories only included trees on public land, we added points where necessary to ensure coverage of all trees visible in the imagery, including trees on private land such as in front yards, backyards, and parking lots. We also visually checked the accuracy of each point and moved it if necessary to ensure that it was on top of the tree trunk location in the image. If we

**Table 1**

Summary of our dataset of annotated NAIP tiles in California. The number of trees listed in the three right-most columns is the number of hand-annotated trees in the study area in each city and in each year. A dash indicates that the area was not annotated for that year. In total, our dataset contains 95,972 tree annotations. The climate zone abbreviations are Inland Empire (IE), Inland Valleys (IV), Interior West (IW), Northern California Coast (NCC), Southern California Coast (SCC), and Southwest Desert (SD). The cities above the line are in Southern California, and below are in Northern California. The cells highlighted in gray indicate the Southern California 2020 portion of the dataset.

City	Zone	Images	Tree annotations		
			2016	2018	2020
Claremont	IE	92	4,856	4,794	4,668
Long Beach	SCC	100	6,470	6,402	5,843
Palm Springs	SD	100	4,431	4,704	4,107
Riverside	IE	90	5,015	4,399	4,082
Santa Monica	SCC	92	5,822	5,829	5,841
Bishop	IW	10	–	–	682
Chico	IV	99	–	8,185	8,162
Eureka	NCC	21	–	–	2,133
<b>Total</b>		<b>1,651</b>	<b>26,594</b>	<b>34,313</b>	<b>35,518</b>

**Table 2**

Summary of dataset subsets used in our experiments.

Subset	Images		Annotations	
	Train	Test	Train	Test
Southern California 2020	426	48	21,861	2,680
Southern California 2016–2018	–	948	–	52,722
Northern California 2018–2020	–	229	–	19,162
<b>Complete dataset</b>	<b>1485</b>	<b>166</b>	<b>87,666</b>	<b>8,759</b>

could not determine the tree trunk location by visual inspection, we put the point on the center of the canopy. This happened sometimes with palm trees, for example, which have thin trunks that are easily confused with their shadow. We compared with imagery at a higher spatial resolution where available to verify tree locations and to ensure that no non-trees (such as shrubs) were included in the annotations.

### 3. Methods

We introduce a novel tree detection method using neural network confidence map regression. Because our annotations are tree location points, we found a confidence map approach (Lempitsky and Zisserman, 2010; Osco et al., 2020) to be the most appropriate for our task. Our overall process for tree detection and localization is illustrated in Fig. 2. The input to the system is a stack of raw and derived rasters. The raster stack is passed through a CNN which outputs a single-channel confidence map. Peaks in the confidence map should correspond to tree locations. We identify tree locations in the confidence map using local peak finding (non-maxima suppression).

Our method is based on the approach of Osco et al. (2020) with several modifications to improve the results. Osco et al. (2020) employ a network using a VGG-16 backbone (Simonyan and Zisserman, 2015) followed by several convolutional blocks with residual connections. They apply peak finding to the confidence map but use a fixed detection threshold without hyperparameter tuning. Our method improves upon the method of Osco et al. (2020) by replacing their network architecture with an attention-based architecture and using hyperparameter tuning to optimize the detection results. In our evaluation comparing our method with several previous methods, we found that our proposed method produced the best results across several metrics.

#### 3.1. Network architecture

We replace the CNN architecture used previously by Osco et al. (2020) with a network based on SFANet (Zhu et al., 2019), an architecture that performed well on several crowd counting benchmarks.



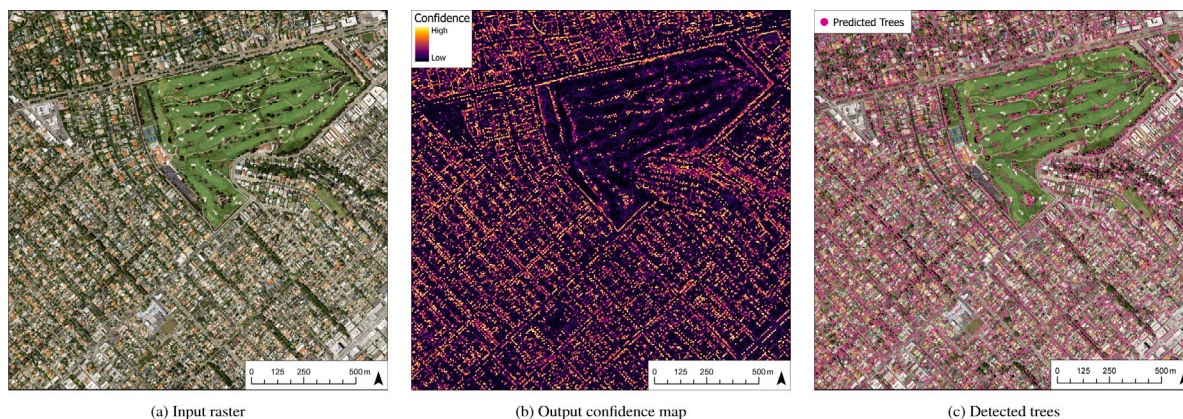


Fig. 2. Example of tree detection using our method on a section of 2020 NAIP imagery in Santa Monica, CA. We process the input raster (a) in a CNN to produce a confidence map (b). We then apply peak finding in the confidence map to produce individual tree detections (c).

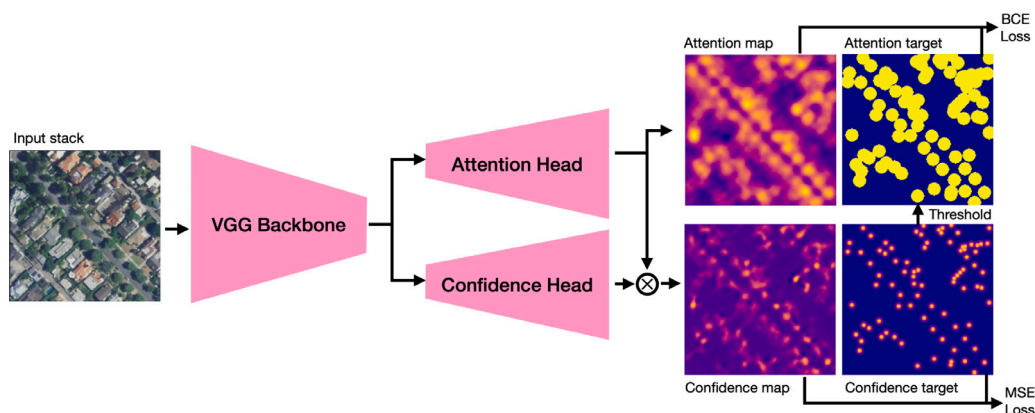


Fig. 3. HR-SFANet network architecture. The input image is encoded through the first five blocks of convolutional and pooling layers from the VGG-16 network. Separate attention and confidence heads upsample and aggregate the outputs of the backbone layers to produce an attention map and a confidence map. These are multiplied together to produce the final confidence map.

Our architecture is illustrated in Fig. 3. SFANet consists of a VGG-16 backbone (Simonyan and Zisserman, 2015), a confidence head, and an attention head. The VGG-16 backbone consists of a series of convolutions and pooling operations to encode the input, while the output heads consist of convolutions, upsampling operations, and skip connections to decode the outputs. We modified the SFANet architecture by adding extra layers to the confidence and attention heads so that they would output at full-resolution rather than half-resolution. We call our network the HR-SFANet. The architecture is illustrated in Fig. 3 and documented in detail in Appendix.

### 3.2. Data preparation

We use NAIP multi-spectral digital number (DN) data as input, which has red (R), green (G), blue (B), and near infrared (N) channels. We also added an NDVI channel (Kriegler et al., 1969; Rouse et al., 1974) derived from the red and near-infrared channels as follows:

$$V = \frac{N - R}{N + R} \tag{1}$$

The NDVI index is a well-known indicator of live green vegetation, and we hypothesized its inclusion would help the network identify trees in the imagery.

Each band of the input is normalized before being processed by the network. The NAIP multispectral DN data is eight-bit and has a range of 0–255. The RGB bands are zero-centered with respect to the ImageNet dataset (Deng et al., 2009) on which the VGG16 backbone network was pre-trained. The *N* band is normalized by subtracting 127.5. The *V* band has a range of −1 to 1 and so we multiply by 127.5 to bring it into a similar range as the other bands after normalization.

### 3.3. Network initialization

The backbone portion of the HR-SFANet network uses a VGG-16 network (Simonyan and Zisserman, 2015) with weights pre-trained on ImageNet (Deng et al., 2009). Since ImageNet contains RGB images, the VGG-16 network is not designed for inputs containing the extra near-IR and NDVI channels present in our imagery. We remedied this by adding two extra input channels to the first layer of the VGG-16 network. The filter weights for these channels are randomly initialized using “Glorot” uniform initialization (Glorot and Bengio, 2010).

### 3.4. Confidence map preparation

During training, each input image tile has a corresponding ground truth confidence map that is used as target for the output of the

network. Following [Osco et al. \(2020\)](#) we form the target confidence map by placing a Gaussian on each tree location and aggregating them with a per-pixel maximum:

$$C(x, y) = \max_i \exp\left(-\frac{(x - x_i)^2 + (y - y_i)^2}{2\sigma^2}\right) \quad (2)$$

where  $(x, y)$  is the location of a pixel in the confidence map and  $(x_i, y_i)$  is the location of the  $i$ th tree in the image. Using the max operator instead of a sum ([Lempitsky and Zisserman, 2010](#)) ensures that objects have distinct peaks in the confidence map, even if they are close to each other.

While we would ideally adapt the Gaussian width parameter  $\sigma$  to each tree based on its size, our point annotations do not contain size information. Therefore we choose a single fixed value for  $\sigma$  for all trees. We tested several settings for  $\sigma$  and found  $\sigma = 1.8$  m to be the best setting in our experiments.

### 3.5. Loss functions

Our primary loss function is the mean squared error (MSE) loss between the predicted and ground truth confidence maps:

$$L_{\text{MSE}} = \frac{1}{N} \sum_{x,y} [C'(x, y) - C(x, y)]^2 \quad (3)$$

where  $C$  and  $C'$  are the ground truth and predicted confidence maps, respectively.

Unlike the CNN used by [Osco et al. \(2020\)](#), our HR-SFANet has an attention head in addition to the confidence head. Following [Zhu et al. \(2019\)](#), we apply the binary cross-entropy (BCE) loss function to the output of the attention head:

$$L_{\text{BCE}} = -\frac{1}{N} \sum_{x,y} [A(x, y) \log(A'(x, y)) + (1 - A(x, y)) \log(1 - A'(x, y))] \quad (4)$$

where  $A$  is a binary mask produced by applying a threshold  $\tau$  to the ground truth confidence map, and  $A'$  is the output of the attention head.

The two loss functions are combined in a weighted sum:

$$L = L_{\text{MSE}} + \alpha L_{\text{BCE}} \quad (5)$$

where  $\alpha$  balances the two loss terms.

Following [Zhu et al. \(2019\)](#), we used  $\tau = 0.001$  and  $\alpha = 0.01$ . We tested different settings of  $\alpha$  but found that the model's performance was not sensitive to the setting of  $\alpha$ .

### 3.6. Peak finding and hyperparameter tuning

During inference, we produce a confidence map from the input raster and then use peak finding to determine predicted tree locations ([Fig. 2](#)). A location in the confidence map is labeled as a peak if it is the local maximum in a region of radius  $2d + 1$ , where  $d$  is the minimum allowable distance between peaks. We also remove local maxima below a threshold, to remove false detections. The threshold can either be an absolute threshold  $t_{\text{abs}}$  or a relative threshold  $t_{\text{rel}}$  (relative to the maximum value in the confidence map).

[Osco et al. \(2020\)](#) proposed to use an absolute threshold of  $t_{\text{abs}} = 0.2$  and a minimum distance of  $d = 3$ . However, we found that we could increase detector performance by finding optimal settings through hyperparameter tuning. For each configuration of the network tested, after training the model, we used hyperparameter tuning to determine the optimal settings for  $d$ ,  $t_{\text{abs}}$ , and  $t_{\text{rel}}$  and to choose whether to use the absolute or relative threshold. We used Optuna hyperparameter optimization framework ([Akiba et al., 2019](#)) and searched for the optimal set of hyperparameters over 200 iterations. We retained the settings with the highest F-score over the validation set. In our experiments, we found that, in contrast to the recommendation of [Osco et al. \(2020\)](#), the relative threshold was always the better choice.

### 3.7. Implementation details

We implemented our HR-SFANet in Python using the Tensorflow ([Abadi et al., 2016](#)) and Keras ([Chollet et al., 2015](#)) libraries. Training and inference were performed on an Nvidia V100 GPU.

During training, we processed batches of  $256 \times 256$  tiles from the training set, using a batch size of eight over 500 epochs. We used the Adam optimizer ([Kingma and Ba, 2015](#)) with learning rate 0.0001,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ , and  $\epsilon = 1 \times 10^{-7}$ . We reserved 10% of the training data for validation and retained the model with best validation loss during training. For data augmentation, we rotated each training patch by  $90^\circ$ ,  $180^\circ$ , and  $270^\circ$  and also horizontally flipped each patch. This increased the size of the training set eight-fold. We experimented with further data augmentation using random image rotation and brightness/color variation but found that it did not lead to an improvement.

### 3.8. Other tree detection methods

We also tested two other well-known tree detection methods in our evaluation: PyCrown ([Zörner et al., 2018](#)) and DeepForest ([Weinstein et al., 2019](#)).

#### 3.8.1. PyCrown

PyCrown ([Zörner et al., 2018](#)) uses LiDAR data as input and applies peak finding on a smoothed canopy height model (CHM) raster to determine tree locations, with some filtering and post-processing applied to eliminate false positives and increase localization accuracy. The method is an open-source re-implementation of [Dalponte and Coomes \(2016\)](#) in Python with some modifications to improve speed and produce better results. Because PyCrown tends to produce many false positives on buildings, we removed any tree detections that landed within OpenStreetMap building extents ([Haklay and Weber, 2008](#)). Since public LiDAR data is not currently available in Palm Springs, CA, we removed those regions from the test set when evaluating PyCrown.

#### 3.8.2. DeepForest

DeepForest ([Weinstein et al., 2019](#)) is a Python package that provides a pre-trained tree detection model based on the RetinaNet object detector ([Lin et al., 2020](#)). The pre-trained model was trained on RGB imagery from sites in the National Ecological Observatory Network (NEON). The method outputs a bounding box for each detected tree.

## 4. Experiments

### 4.1. Evaluation metrics

We evaluate our models on five standard metrics: Average Precision (AP), precision, recall, F-score, and Root Mean Square Error (RMSE). Given the count of true positives (TP), false positives (FP), and false negatives (FN), precision, recall, and F-score are calculated as follows:

$$\text{precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \quad (6)$$

$$\text{recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (7)$$

$$\text{F-score} = 2 \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} \quad (8)$$

True positives, false positives, and false negatives are determined by matching predicted tree locations and ground truth tree locations. We find the minimum weight matching, which minimizes the sum of distances between matched trees, and ensures that each ground truth tree and each predicted tree has at most one match. We remove matches whose distance is over a threshold; in our evaluation we used a threshold of six meters. Example results from the matching procedure can be seen in [Fig. 4](#).

RMSE is calculated from the distance between true positive predictions and their corresponding ground truth trees. Let  $\mathcal{T}$  be the set

**Table 3**  
Comparison of tree detection methods on Southern California 2020 test set.

Method	AP	Precision	Recall	F-Score	RMSE [m]
PyCrown (Zörner et al., 2018)	–	0.401	0.661	0.499	2.709
DeepForest (Weinstein et al., 2019)	0.387	0.735	0.294	0.420	2.719
+ fine-tuning	0.701	0.707	0.713	0.710	2.431
+ hyperparameter tuning	0.701	0.782	0.683	0.729	2.413
Osco et al. (2020)	0.660	<b>0.803</b>	0.476	0.598	2.270
+ hyperparameter tuning	0.660	0.764	0.706	0.734	2.263
HR-SFANet (ours)	<b>0.705</b>	0.736	<b>0.733</b>	<b>0.735</b>	<b>2.157</b>

**Table 4**  
Results of our method when extrapolating to different climate zones and years. The top portion consists of the same Southern California cities in the 2020 training set, but earlier image capture years, while the bottom portion consists of Northern California cities distinct from the training set region.

City	Year	Zone	AP	Precision	Recall	F-Score	RMSE [m]
Santa Monica	2016	SCC	0.665	0.734	0.693	0.713	2.132
Santa Monica	2018	SCC	0.707	0.732	0.730	0.731	2.188
Long Beach	2016	SCC	0.694	0.799	0.668	0.728	1.963
Long Beach	2018	SCC	0.720	0.856	0.639	0.732	1.830
Claremont	2016	IE	0.654	0.739	0.668	0.701	2.075
Claremont	2018	IE	0.642	0.708	0.668	0.687	2.435
Riverside	2016	IE	0.723	0.819	0.671	0.737	1.905
Riverside	2018	IE	0.590	0.686	0.622	0.652	2.599
Palm Springs	2016	SD	0.624	0.700	0.647	0.672	1.953
Palm Springs	2018	SD	0.645	0.743	0.620	0.676	1.804
Chico	2018	IV	0.708	0.748	0.688	0.716	2.179
Chico	2020	IV	0.701	0.733	0.689	0.710	2.338
Eureka	2020	NCC	0.557	0.744	0.509	0.604	2.423
Bishop	2020	IW	0.687	0.723	0.694	0.708	2.203

of accepted matches and let  $(x_i, y_i)$  and  $(\hat{x}_i, \hat{y}_i)$  be the locations of corresponding predicted and ground truth trees, respectively.

$$RMSE = \sqrt{\frac{1}{|\mathcal{T}|} \sum_{i \in \mathcal{T}} (x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2} \tag{9}$$

To calculate AP, we calculate precision and recall over a range of absolute thresholds for accepting predicted trees according to their confidence value (the value of the confidence map at the tree’s location). Let  $P_n$  and  $R_n$  be the precision and recall observed at the  $n$ th threshold. The AP is calculated as:

$$AP = \sum_n (R_n - R_{n-1}) P_n \tag{10}$$

AP is useful because it summarizes the performance of the detector across a range of detection thresholds, whereas the other metrics describe the performance of the detector at a single detection threshold optimized to balance both precision and recall.

To calculate evaluation metrics for DeepForest, we considered each ground truth point contained within a predicted bounding box as a possible true positive. We calculated an optimal matching to determine the assignment between ground truth points and bounding boxes which maximized the number of true positives. To calculate RMSE for DeepForest we used the center points of the bounding boxes for the predicted tree locations.

## 4.2. Comparison of methods on Southern California test set

We tested each method on the Southern California 2020 portion of the dataset. The results are summarized in Table 3.

### 4.2.1. Our method

As our method is based on Osco et al. (2020) we first tested their method on our dataset. Since their code is not available, we created our own implementation based on the description in the paper and trained the network from Osco et al. (2020) in the same manner as our method. Using their recommended settings for peak finding, the Osco

et al. method achieves an AP of 0.660, F-Score of 0.598, and RMSE of 2.270 m. After applying our proposed method of hyperparameter tuning to find the optimal peak finding settings, F-Score increases to 0.734, and RMSE decreases to 2.263 m. When we replace their network with our HR-SFANet, AP increases to 0.705, F-Score increases to 0.735, and RMSE decreases to 2.157 m. This experiment validates the importance of our proposed network design and hyperparameter tuning approach to improve the results.

### 4.2.2. PyCrown

PyCrown only achieved an F-score of 0.499 and RMSE of 2.709 m on our test set. We could not calculate AP for PyCrown because there is no detection threshold to vary in this method. The precision for PyCrown was exceptionally low (0.401), indicating that it output more false positives than other methods. This could be due to a single tree having multiple local peaks in its canopy height map, leading to false detections.

### 4.2.3. DeepForest

The pre-trained DeepForest model achieved a high precision (0.735) but low recall (0.294), indicating that it was unable to detect many trees that other methods could detect. This led to an F-Score of 0.420 and AP of 0.387. This poor performance is likely due to the fact that the NEON imagery on which DeepForest was trained is higher resolution than our NAIP imagery, and mostly contains wilderness areas whereas our imagery contains urban environments.

To improve the results of DeepForest, we fine-tuned the pre-trained model on our training set. Because our data does not include bounding box annotations, we placed a fixed size bounding box around each ground truth point. We experimented with different box sizes and chose the size that led to the best results on the test set, which was a box with width and height of 9.6 m (16 pixels). We fine-tuned the model for 15 epochs; training past this point did not yield any further improvement. After fine-tuning, the AP increased substantially to 0.701, and the F-Score increased to 0.710 using the default detection threshold setting. After hyper-parameter tuning, the F-Score improved to 0.729 with an RMSE of 2.413. This experiment indicates that a bounding box object detector can be trained to effectively detect trees using point annotations; however, the performance is below what we achieved with our confidence map approach.

### 4.2.4. Inference speed

We calculated the average processing time for each method to detect trees in a 256 × 256 raster. PyCrown is the slowest method, taking over five seconds per image. The CNN-based methods are much faster, and ours is the fastest. DeepForest takes 71 ms per image, Osco et al. (2020) takes 44 ms per image, and our method takes 21 ms per image.

## 4.3. Extrapolation to other years and climate zones

When trained and tested on 2020 imagery from Southern California, our method achieved the best results among methods tested. We then evaluated the performance of our method on pre-2020 Southern California imagery and on imagery from three cities in Northern California, each in a different climate zone. None of the images in these experiments were present in the training or validation data. Because





Fig. 4. Example results from our tree detection method on images from Southern California 2020 test set. The method is able to detect and accurately localize most of the trees in the images. Some areas contain examples of missed detections of trees with small canopies and shrubs being confused with trees.

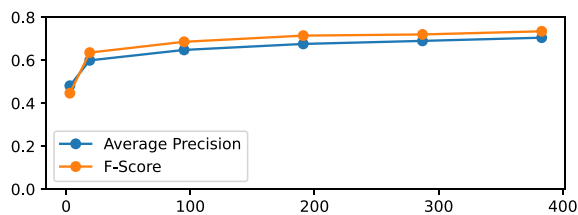


Fig. 5. Analysis of test set performance with increasing training set size. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

these test sets represent different climate zones, varying tree density, and in some cases an earlier image capture year than the training set (see Table 1), they serve to evaluate our method’s ability to generalize to data different from the training set. The results are summarized in Table 4. Overall, the precision in these areas was similar to the precision observed in the 2020 Southern California test set, ranging from 0.856 to 0.686, but recall was lower, ranging from 0.730 to 0.509. The F-score ranged from 0.737 to 0.604.

#### 4.4. Qualitative analysis

Fig. 4 shows example results from our method on selected regions from our test sets. In general, we noticed three typical sources of errors: missed detections of trees with small canopies; confusion of different types of plant, such as shrubs, with trees; and over- or under-estimation of trees in areas of dense canopy, which sometimes have strong shadows. Another challenge is the variation in viewing angle, illumination, and canopy size seen in the NAIP imagery, since the time and date of capture varied across regions and from year to year. These sources of errors would be mitigated by adding more training data or fine-tuning on specific regions, and using imagery with higher spatial and spectral resolution.

Table 5

Test set results for our HR-SFANet model trained on 90% of the entire dataset, covering all years and climate zones. The test set is the remaining 10% of the dataset.

AP	Precision	Recall	F-Score	RMSE [m]
0.727	0.765	0.720	0.741	1.936

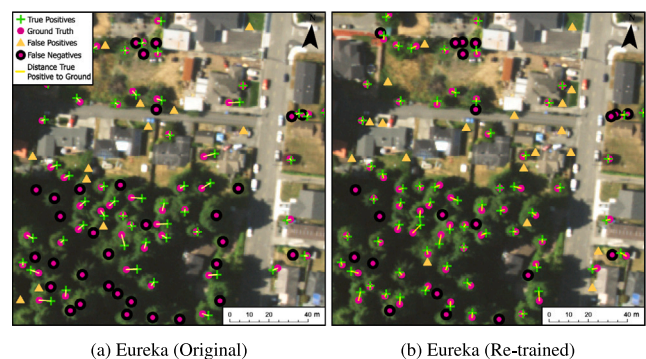


Fig. 6. Results in Eureka (a) before and (b) after re-training model with data from all areas and image years. After re-training, the model detects more trees in the dense and shadowed areas.

#### 4.5. Effect of reducing training set size

Deep learning methods are well-known for requiring large amounts of training data to be effective; however, manual annotation of images is time-intensive and tedious work. We used the Southern California 2020 portion of our dataset to investigate the relationship between training set size and test set performance for our task. We subsampled the training set to various amounts, trained a separate network with



Fig. 7. Comparison of a manual tree inventory with our tree detection results on three areas from Torrance, CA. Our tree detector is able to find trees on both public and private land, leading to a more complete inventory. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

each subset, and calculated the performance metrics for each network using the complete test set. The results are shown in Fig. 5. With 1% of the data (3 images), the F-Score is reduced to 0.447 and AP is 0.481. With 5% of the data (19 images), F-Score improves to 0.636 and AP to 0.600. From there, the metrics more slowly increase with increasing training set size. This experiment shows that there are diminishing gains in adding more training data.

#### 4.6. Effect of increasing training set size

The lowest recall (0.509) was observed in Eureka, a city in Northern California which contains dense natural forest areas within its urban boundary (Fig. 6). To evaluate the potential to address sources of error and improve performance by adding more training data, we re-trained the model using a random 90% subset of the entire dataset (1,485 images), not just the Southern California 2020 portion as in the original model. The results of the original and re-trained model in Eureka are compared in Fig. 6, which illustrates how the re-trained model detects trees in the dense and shadowed region that were missed by the original model. Note that the test image shown in Fig. 6 was not in the training set of either model.

Quantitative results for the re-trained model are shown in Table 5. The re-trained model improved AP to 0.727, F-score to 0.741, and RMSE to 1.936 m. However, note that these numbers are not directly comparable to the results for the models trained and tested on the Southern California 2020 subset (Table 3), because the test sets for the two evaluations are not the same.

## 5. Application

Because the HR-SFANet network in our method is fully convolutional (Shelhamer et al., 2017), it can process any size raster as input. However, a large raster will need to be processed in tiles to stay within the memory limitations of the machine. The naïve approach of processing each tile without overlap between the tiles results in disagreement at the edges of the confidence maps, leading to inaccurate, missing, or duplicated tree detections. To avoid these artifacts when processing a large raster, we divide the raster into an overlapping grid of tiles. After processing by the network, the overlap region is discarded to avoid edge artifacts. We used a tile size of  $2112 \times 2112$  pixels with an overlap of 32 pixels. We similarly apply local peak finding on an overlapping grid when processing a large raster; for peak finding we used a tile size of  $256 \times 256$  and an overlap of 32 pixels.

Using our method with this tiled inference approach, we processed NAIP 2020 imagery of all of California's urban areas. Because our focus is on creating tree inventories for cities, we excluded any tree

detections outside of urban boundaries as determined by the California Department of Water Resources (Land IQ, 2017)

Our automatic tree detection approach can fill deficiencies in tree inventories that typically only cover street trees on public rights-of-way areas and thus miss a large proportion of the trees in a city's urban forest. For example, Fig. 7 compares a manual tree inventory from Torrance, CA (Love et al., 2022) with our automatic tree detection results. The blue points are from the existing urban forest inventory, and the purple points are trees detected with our method. Our tree detector is able to detect trees in private spaces that are not labeled in the public inventory.

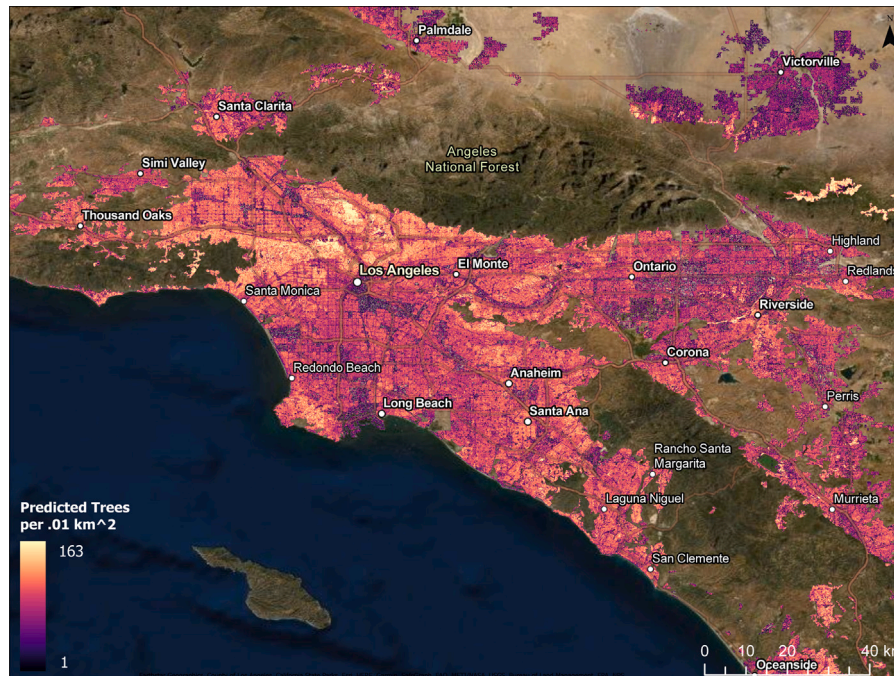
Our method can also support larger-scale analyses of tree density across the state. For example, Fig. 8 shows a tree density map for urban areas in Southern California. To create this map, we calculated the number of detected trees per grid cell over a  $100 \text{ m} \times 100 \text{ m}$  grid.

## 6. Conclusions and future work

Creating good data and efficient systems for the management of our urban forests will be essential as our climate changes, and the ecosystem services, such as controlling microclimate, of our urban forests become increasingly important (McPherson et al., 2016). Currently, inventories of a city's urban forest are completed manually, and repeated on a routine basis (Nielsen et al., 2014). This time-consuming process costs cities significant time and financial resources, and typically only accounts for the trees within the publicly managed urban forest. Inventories are used by cities to manage their urban forest, and used by cities and researchers to estimate the ecosystem services given by their forests (McPherson et al., 2016; Love et al., 2022). Trees on privately managed land, like public trees, contribute to the ecosystem services an urban forest provides to its residents, but there are few ways to account for those services in typical tree inventories and resulting research. Our research allows for trees in the entire urban area to be counted, creating more accurate total tree estimates for cities and neighborhoods. Because our counts are spatially explicit, managers can use them to identify areas with relatively few urban trees, and attempt to address why that is the case and target that area for planting. Our data can be used to compare socioeconomic trends with tree counts, find trends in where there are more publicly or privately managed trees, and identify areas of tree inequity in urban areas. By comparing tree inventories across years, we could analyze the progress of the urban forest over time and better understand the effect of environmental conditions such as drought.

Future work could include exploring the use of imagery with higher spatial and spectral resolution, to help the detector separate out individual tree canopies in dense stands of trees, detect small trees,





**Fig. 8.** Visualization of tree count per  $.01 \text{ km}^2$  in the Los Angeles area. Each  $.01 \text{ km}^2$  is visualized using the total number of trees we detected in that grid after applying the Tree Detector to 2020 NAIP imagery. The basemap shown is the World Imagery Basemap accessed through ArcGIS Pro (ESRI, 2023).

and distinguish trees from other visually similar plants. It also will be important to consider how to easily enable end users to fine-tune the model for their particular region of interest.

Furthermore, while our method accurately detects and localizes trees from aerial imagery, on-the-ground urban tree inventories provide much more information than our current method can provide, such as: coarse-grained classification of the tree (e.g., identifying deciduous and coniferous trees); fine-grained classification of the tree genus or species; estimation of canopy size; health status of the tree; and monitoring of silvicultural treatments. The ability to produce such information would be highly valuable for improving an urban tree inventory. Future work lies in incorporating techniques such as image-based classification (Beery et al., 2022) and canopy size estimation using instance segmentation (Sun et al., 2022) to automatically determine these important tree properties from imagery.

#### CRedit authorship contribution statement

**Jonathan Ventura:** Conceptualization, Investigation, Methodology, Software, Validation, Writing – original draft, Writing – review & editing. **Camille Pawlak:** Conceptualization, Data curation, Validation, Visualization, Writing – original draft, Writing – review & editing. **Milo Honsberger:** Data curation. **Cameron Gonsalves:** Data curation. **Julian Rice:** Data curation, Software. **Natalie L.R. Love:** Conceptualization. **Skyler Han:** Software. **Viet Nguyen:** Software. **Keilana Sugano:** Data curation. **Jacqueline Doremus:** Supervision. **G. Andrew Fricker:** Conceptualization, Methodology, Writing – original draft, Writing – review & editing. **Jenn Yost:** Conceptualization, Funding acquisition, Project administration, Writing – review & editing. **Matt Ritter:** Conceptualization, Funding acquisition, Project administration, Writing – review & editing.

#### Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

#### Data availability

The code is available at <https://github.com/jonathanventura/urban-tree-detection> and the data is available at <https://github.com/jonathanventura/urban-tree-detection-data>. Tree detection results are available at OSF:4S859 (Ventura et al., 2024) and as an interactive map at <https://ufe1.calpoly.edu/>.

#### Acknowledgments

This project was funded by CAL FIRE (award number: 8GB18415) the US Forest Service (award number: 21-CS-11052021-201), and an incubation grant from the Data Science Strategic Research Initiative at California Polytechnic State University. Thanks to Allan Hollander, Jim Thorne, Russ White, Ronny Hänsch, and the anonymous reviewers for their comments on the manuscript.

#### Appendix. HR-SFANet architecture

Here we provide pseudo-code to exactly describe the architecture of our HR-SFANet network, which is a modification of the original SFANet architecture (Zhu et al., 2019). The pseudocode makes use of the following standard neural network operations:

- $\text{CONV}(n, k, x)$ : 2D convolution on input  $x$  with a filter size of  $k \times k$  and  $n$  output channels.
- $\text{ReLU}(x) = \max(0, x)$ : Element-wise rectified linear unit.
- $\text{SIGMOID}(x) = 1/(1 + \exp(-x))$ : Element-wise sigmoid function.
- $\text{BATCHNORM}(x)$ : Batch normalization (Ioffe and Szegedy, 2015).
- $\text{MAXPOOL}(x)$ :  $2 \times 2$  max pooling with a stride of 2.
- $\text{CONCATENATE}(x, y)$ : Channel-wise concatenation of  $x$  and  $y$ .
- $\text{UPSAMPLE}(x)$ :  $2 \times$  upsampling of input  $x$  using bilinear interpolation.

**Algorithm 1** High-Resolution SFANet

---

```

function HR-SFANet(input)
  conv12, conv22, conv33, conv43, conv53 ← VGG16BACKBONE(input)
  attention_map ← DECODER(conv12, conv22, conv33, conv43, conv53)
  confidence_map ← SIGMOID(BATCHNORM(CONV(1, 1, attention_map)))
  confidence_map ← DECODER(conv12, conv22, conv33, conv43, conv53)
  confidence_map ← CONV(1, 1, attention_map * confidence_map)
  return confidence_map, attention_map
end function

```

---

**Algorithm 2** VGG-16 Backbone

---

```

function VGG16BACKBONE(input)
  conv11 ← RELU(BATCHNORM(CONV(64, 3, input)))                                ▷ Block 1
  conv12 ← RELU(BATCHNORM(CONV(64, 3, conv11)))
  conv21 ← RELU(BATCHNORM(CONV(128, 3, MAXPOOL(conv21))))                    ▷ Block 2
  conv22 ← RELU(BATCHNORM(CONV(128, 3, conv21)))
  conv31 ← RELU(BATCHNORM(CONV(256, 3, MAXPOOL(conv22))))                    ▷ Block 3
  conv32 ← RELU(BATCHNORM(CONV(256, 3, conv31)))
  conv33 ← RELU(BATCHNORM(CONV(256, 3, conv32)))
  conv41 ← RELU(BATCHNORM(CONV(512, 3, MAXPOOL(conv33))))                    ▷ Block 4
  conv42 ← RELU(BATCHNORM(CONV(512, 3, conv41)))
  conv43 ← RELU(BATCHNORM(CONV(512, 3, conv42)))
  conv51 ← RELU(BATCHNORM(CONV(512, 3, MAXPOOL(conv43))))                    ▷ Block 5
  conv52 ← RELU(BATCHNORM(CONV(512, 3, conv51)))
  conv53 ← RELU(BATCHNORM(CONV(512, 3, conv52)))
  return conv12, conv22, conv33, conv43, conv53
end function

```

---

**Algorithm 3** Decoder

---

```

function DECODER(conv11, conv22, conv33, conv43, conv53)
  x ← CONCATENATE(UPSAMPLE(conv53), conv43)                                ▷ Block 1
  x ← RELU(BATCHNORM(CONV(256, 1, x)))
  x ← RELU(BATCHNORM(CONV(256, 3, x)))
  x ← CONCATENATE(UPSAMPLE(x), conv33)                                    ▷ Block 2
  x ← RELU(BATCHNORM(CONV(128, 1, x)))
  x ← RELU(BATCHNORM(CONV(128, 3, x)))
  x ← CONCATENATE(UPSAMPLE(x), conv22)                                    ▷ Block 3
  x ← RELU(BATCHNORM(CONV(64, 1, x)))
  x ← RELU(BATCHNORM(CONV(64, 3, x)))
  x ← RELU(BATCHNORM(CONV(32, 3, x)))
  x ← CONCATENATE(UPSAMPLE(x), conv12)                                    ▷ Block 4
  x ← RELU(BATCHNORM(CONV(32, 1, x)))
  x ← RELU(BATCHNORM(CONV(32, 3, x)))
  x ← RELU(BATCHNORM(CONV(32, 3, x)))
  return x
end function

```

---

**References**

- Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., Kudlur, M., Levenberg, J., Monga, R., Moore, S., Murray, D.G., Steiner, B., Tucker, P., Vasudevan, V., Warden, P., Wicke, M., Yu, Y., Zheng, X., 2016. TensorFlow: A system for large-scale machine learning. In: Proceedings of the 12th USENIX Conference on Operating Systems Design and Implementation. OSDI '16, USENIX Association, USA, pp. 265–283.
- Akiba, T., Sano, S., Yanase, T., Ohta, T., Koyama, M., 2019. Optuna: A next-generation hyperparameter optimization framework. In: Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. KDD '19, Association for Computing Machinery, New York, NY, USA, pp. 2623–2631. <http://dx.doi.org/10.1145/3292500.3330701>.
- Alonzo, M., Roth, K., Roberts, D., 2013. Identifying Santa Barbara's urban tree species from AVIRIS imagery using canonical discriminant analysis. Remote Sens. Lett. 4 (5), 513–521. <http://dx.doi.org/10.1080/2150704X.2013.764027>.
- Ayrey, E., Hayes, D.J., 2018. The use of three-dimensional convolutional neural networks to interpret LiDAR for forest inventory. Remote Sens. 10 (4), <http://dx.doi.org/10.3390/rs10040649>.
- Ball, J.G., Hickman, S.H., Jackson, T.D., Koay, X.J., Hirst, J., Jay, W., Archer, M., Aubry-Kientz, M., Vincent, G., Coomes, D.A., 2023. Accurate delineation of individual tree crowns in tropical forests from aerial RGB imagery using Mask R-CNN. Remote Sens. Ecol. Conserv. 9 (5), 641–655.
- Beery, S., Wu, G., Edwards, T., Pavetic, F., Majewski, B., Mukherjee, S., Chan, S., Morgan, J., Rathod, V., Huang, J., 2022. The Auto Arborist Dataset: A large-scale benchmark for multiview urban forest monitoring under domain shift. In: 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition. CVPR, pp. 21262–21275. <http://dx.doi.org/10.1109/CVPR52688.2022.02061>.
- Belouï, M., Heinzmann, L., Rehush, N., Gessler, A., Griess, V.C., 2023. Individual tree-crown detection and species identification in heterogeneous forests using aerial RGB imagery and deep learning. Remote Sens. 15 (5), 1463.
- Bosch, M., 2020. Detectree: Tree detection from aerial imagery in python. J. Open Source Softw. 5 (50), 2172. <http://dx.doi.org/10.21105/joss.02172>.
- Brandt, M., Tucker, C.J., Kariyaa, A., Rasmussen, K., Abel, C., Small, J., Chave, J., Rasmussen, L.V., Hiernaux, P., Diouf, A.A., et al., 2020. An unexpectedly large count of trees in the West African Sahara and Sahel. Nature 587 (7832), 78–82. <http://dx.doi.org/10.1038/s41586-020-2824-5>.

- Branson, S., Wegner, J.D., Hall, D., Lang, N., Schindler, K., Perona, P., 2018. From Google Maps to a fine-grained catalog of street trees. *ISPRS J. Photogramm. Remote Sens.* 135, 13–30. <http://dx.doi.org/10.1016/j.isprsjprs.2017.11.008>.
- Chen, X., Jiang, K., Zhu, Y., Wang, X., Yun, T., 2021. Individual tree crown segmentation directly from UAV-Borne LiDAR data using the PointNet of deep learning. *Forests* 12 (2), <http://dx.doi.org/10.3390/f12020131>.
- Chen, G., Shang, Y., 2022. Transformer for tree counting in Aerial images. *Remote Sens.* 14 (3), <http://dx.doi.org/10.3390/rs14030476>.
- Chen, G., Zakhor, A., 2009. 2D tree detection in large urban landscapes using aerial LiDAR data. In: 2009 16th IEEE International Conference on Image Processing. ICIP, pp. 1693–1696. <http://dx.doi.org/10.1109/ICIP.2009.5413699>.
- Chollet, F., et al., 2015. Keras. <https://github.com/fchollet/keras>.
- Clarke, L.W., Jenerette, G.D., Davila, A., 2013. The luxury of vegetation and the legacy of tree biodiversity in Los Angeles, CA. *Landsc. Urban Plan.* 116, 48–59. <http://dx.doi.org/10.1016/j.landurbplan.2013.04.006>.
- Dalponte, M., Coomes, D.A., 2016. Tree-centric mapping of forest carbon density from airborne laser scanning and hyperspectral data. *Methods Ecol. Evol.* 7 (10), 1236–1245. <http://dx.doi.org/10.1111/2041-210X.12575>.
- Das, S., Sun, Q., Zhou, H., 2022. GeoAI to implement an individual tree inventory: Framework and application of heat mitigation. *Urban For. Urban Green.* 74, 127634. <http://dx.doi.org/10.1016/j.ufug.2022.127634>.
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., Fei-Fei, L., 2009. ImageNet: A large-scale hierarchical image database. In: 2009 IEEE Conference on Computer Vision and Pattern Recognition. pp. 248–255. <http://dx.doi.org/10.1109/CVPR.2009.5206848>.
- Dwyer, J.F., McPherson, E.G., Schroeder, H.W., Rowntree, R.A., 1992. Assessing the benefits and costs of the urban forest. *J. Arboric.* 18, 227.
- ESRI, 2023. ArcGIS Pro.
- Freudenberg, M., Magdon, P., Nölke, N., 2022. Individual tree crown delineation in high-resolution remote sensing images based on U-Net. *Neural Comput. Appl.* 34 (24), 22197–22207. <http://dx.doi.org/10.1007/s00521-022-07640-4>.
- Glorot, X., Bengio, Y., 2010. Understanding the difficulty of training deep feedforward neural networks. In: Teh, Y.W., Titterton, M. (Eds.), *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*. In: *Proceedings of Machine Learning Research*, Vol. 9, PMLR, Chia Laguna Resort, Sardinia, Italy, pp. 249–256.
- Haklay, M., Weber, P., 2008. OpenStreetMap: User-generated street maps. *IEEE Pervasive Comput.* 7 (4), 12–18. <http://dx.doi.org/10.1109/MPRV.2008.80>.
- Ioffe, S., Szegedy, C., 2015. Batch Normalization: Accelerating deep network training by reducing internal covariate shift. In: *International Conference on Machine Learning*. PMLR, pp. 448–456.
- Jensen, R.R., Hardin, P.J., Hardin, A.J., 2012. Classification of urban tree species using hyperspectral imagery. *Geocarto Int.* 27 (5), 443–458. <http://dx.doi.org/10.1080/10106049.2011.638989>.
- Kingma, D.P., Ba, J., 2015. Adam: A method for stochastic optimization. In: *International Conference on Learning Representations*. <http://dx.doi.org/10.48550/arXiv.1412.6980>.
- Kriegler, F., Malila, W., Nalepka, R., Richardson, W., 1969. Preprocessing transformations and their effects on multispectral recognition. In: *Proceedings of the 6th International Symposium on Remote Sensing of Environment*. pp. 97–131.
- Land IQ, L., 2017. Land Use - 2014 - Land IQ [ds2677]. <https://map.dfg.ca.gov/metadata/ds2677.html>.
- Lempitsky, V., Zisserman, A., 2010. Learning to count objects in images. In: Lafferty, J., Williams, C., Shawe-Taylor, J., Zemel, R., Culotta, A. (Eds.), *In: Advances in Neural Information Processing Systems*, Vol. 23, Curran Associates, Inc..
- Lin, T.-Y., Goyal, P., Girshick, R., He, K., Dollár, P., 2020. Focal loss for dense object detection. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence*. Vol. 42, pp. 318–327. <http://dx.doi.org/10.1109/TPAMI.2018.2858826>.
- Liu, T., Im, J., Quackenbush, L.J., 2015. A novel transferable individual tree crown delineation model based on Fishing Net Dragging and boundary classification. *ISPRS J. Photogramm. Remote Sens.* 110, 34–47. <http://dx.doi.org/10.1016/j.isprsjprs.2015.10.002>.
- Livesley, S.J., McPherson, E.G., Calfapietra, C., 2016. The urban forest and ecosystem services: Impacts on urban water, heat, and pollution cycles at the Tree, Street, and City Scale. *J. Environ. Qual.* 45 (1), 119–124. <http://dx.doi.org/10.2134/jeq2015.11.0567>.
- Love, N.L., Nguyen, V., Pawlak, C., Pineda, A., Reimer, J.L., Yost, J.M., Fricker, G.A., Ventura, J.D., Doremus, J.M., Crow, T., Ritter, M.K., 2022. Diversity and structure in California's urban forest: What over six million data points tell us about one of the world's largest urban forests. *Urban For. Urban Green.* 74, 127679. <http://dx.doi.org/10.1016/j.ufug.2022.127679>.
- Ma, B., Hauer, R.J., Wei, H., Koeser, A.K., Peterson, W., Simons, K., Timilsina, N., Werner, L.P., Xu, C., 2020. An assessment of street tree diversity: Findings and implications in the United States. *Urban For. Urban Green.* 56, 126826. <http://dx.doi.org/10.1016/j.ufug.2020.126826>.
- Martins, G.B., La Rosa, L.E.C., Happ, P.N., Filho, L.C.T.C., Santos, C.J.F., Feitosa, R.Q., Ferreira, M.P., 2021. Deep learning-based tree species mapping in a highly diverse tropical urban setting. *Urban For. Urban Green.* 64, 127241. <http://dx.doi.org/10.1016/j.ufug.2021.127241>.
- McPherson, E.G., 2010. Selecting reference cities for i-Tree Streets. *Arboric. Urban For.* 36 (5), 230–240. <http://dx.doi.org/10.48044/jauf.2010.031>.
- McPherson, E.G., Simpson, J.R., 2002. A comparison of municipal forest benefits and costs in Modesto and Santa Monica, California, USA. *Urban For. Urban Green.* 1 (2), 61–74. <http://dx.doi.org/10.1078/1618-8667-00007>.
- McPherson, E.G., van Doorn, N., de Goede, J., 2016. Structure, function and value of street trees in California, USA. *Urban For. Urban Green.* 17, 104–115. <http://dx.doi.org/10.1016/j.ufug.2016.03.013>.
- McPherson, E.G., Xiao, Q., van Doorn, N.S., de Goede, J., Bjorkman, J., Hollander, A., Boynton, R.M., Quinn, J.F., Thorne, J.H., 2017. The structure, function and value of urban forests in California communities. *Urban For. Urban Green.* 28, 43–53. <http://dx.doi.org/10.1016/j.ufug.2017.09.013>.
- Münzinger, M., Prechtel, N., Behnisch, M., 2022. Mapping the urban forest in detail: From LiDAR point clouds to 3D tree models. *Urban For. Urban Green.* 74, 127637. <http://dx.doi.org/10.1016/j.ufug.2022.127637>.
- Nielsen, A.B., Östberg, J., Delshammar, T., et al., 2014. Review of urban tree inventory methods used to collect data at single-tree level. *Arboric. Urban For.* 40 (2), 96–111. <http://dx.doi.org/10.48044/jauf.2014.011>.
- Oscro, L.P., dos Santos de Arruda, M., Marcato Junior, J., da Silva, N.B., Ramos, A.P.M., Moryia, É.A.S., Imai, N.N., Pereira, D.R., Creste, J.E., Matsubara, E.T., Li, J., Gonçalves, W.N., 2020. A convolutional neural network approach for counting and geolocating citrus-trees in UAV multispectral imagery. *ISPRS J. Photogramm. Remote Sens.* 160, 97–106. <http://dx.doi.org/10.1016/j.isprsjprs.2019.12.010>.
- Rouse, J., Haas, R., J.A. Scheel, J., Deering, D., 1974. *Monitoring vegetation systems in the Great Plains with ERTS*. In: *Proceedings, 3rd Earth Resource Technology Satellite (ERTS) Symposium*. Vol. 1, pp. 48–62.
- Roussel, J.-R., Auty, D., Coops, N.C., Tompalski, P., Goodbody, T.R., Meador, A.S., Bourdon, J.-F., De Boissieu, F., Achim, A., 2020. lidR: An R package for analysis of Airborne Laser Scanning (ALS) data. *Remote Sens. Environ.* 251, 112061. <http://dx.doi.org/10.1016/j.rse.2020.112061>.
- Santos, A.A.d., Marcato Junior, J., Araújo, M.S., Di Martini, D.R., Tetila, E.C., Siqueira, H.L., Aoki, C., Eltner, A., Matsubara, E.T., Pistori, H., et al., 2019. Assessment of CNN-based methods for individual tree detection on images captured by RGB cameras attached to UAVs. *Sensors* 19 (16), 3595. <http://dx.doi.org/10.3390/s19163595>.
- Shang, X., Chisholm, L.A., 2014. Classification of Australian native forest species using hyperspectral remote sensing and machine-learning classification algorithms. *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.* 7 (6), 2481–2489. <http://dx.doi.org/10.1109/JSTARS.2013.2282166>.
- Shelhamer, E., Long, J., Darrell, T., 2017. Fully convolutional networks for semantic segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.* 39 (4), 640–651. <http://dx.doi.org/10.1109/TPAMI.2016.2572683>.
- Silva, C.A., Hudak, A.T., Vierling, L.A., Loudermilk, E.L., O'Brien, J.J., Hiers, J.K., Jack, S.B., Gonzalez-Benecke, C., Lee, H., Falkowski, M.J., Khosravipour, A., 2016. Imputation of individual longleaf pine (*Pinus palustris* Mill.) Tree attributes from field and LiDAR data. *Can. J. Remote Sens.* 42 (5), 554–573. <http://dx.doi.org/10.1080/07038992.2016.1196582>.
- Simonyan, K., Zisserman, A., 2015. Very deep convolutional networks for large-scale image recognition. In: Bengio, Y., LeCun, Y. (Eds.), *3rd International Conference on Learning Representations, ICLR*. <http://dx.doi.org/10.48550/arXiv.1409.1556>.
- Sun, Y., Li, Z., He, H., Guo, L., Zhang, X., Xin, Q., 2022. Counting trees in a subtropical mega city using the instance segmentation method. *Int. J. Appl. Earth Obs. Geoinf.* 106, 102662. <http://dx.doi.org/10.1016/j.jag.2021.102662>.
- U.S. Census Bureau, 2020. *Decennial census population*.
- Ventura, J., Pawlak, C., Honsberger, M., Gonsalves, C., Rice, J., Love, N., Han, S., Nguyen, V., Sugano, K., Doremus, J., Fricker, G.A., Yost, J., Ritter, M., 2024. Results: Individual tree detection in large-scale urban environments using high-resolution multispectral imagery. <http://dx.doi.org/10.17605/OSF.IO/4S859>, URL <https://osf.io/4s859>.
- Wegner, J.D., Branson, S., Hall, D., Schindler, K., Perona, P., 2016. Cataloging public objects using aerial and street-level images - urban trees. In: 2016 IEEE Conference on Computer Vision and Pattern Recognition. CVPR, pp. 6014–6023. <http://dx.doi.org/10.1109/CVPR.2016.647>.
- Weinstein, B.G., Marconi, S., Bohlman, S., Zare, A., White, E., 2019. Individual tree-crown detection in RGB imagery using semi-supervised deep learning neural networks. *Remote Sens.* 11 (11), 1309. <http://dx.doi.org/10.3390/rs11111309>.
- Wu, B., Yu, B., Wu, Q., Huang, Y., Chen, Z., Wu, J., 2016. Individual tree crown delineation using localized contour tree method and airborne LiDAR data in coniferous forests. *Int. J. Appl. Earth Obs. Geoinf.* 52, 82–94. <http://dx.doi.org/10.1016/j.jag.2016.06.003>.



- Xiao, Q., Ustin, S.L., McPherson, E.G., 2004. Using AVIRIS data and multiple-masking techniques to map urban forest tree species. *Int. J. Remote Sens.* 25 (24), 5637–5654. <http://dx.doi.org/10.1080/01431160412331291224>.
- Xu, W., Deng, S., Liang, D., Cheng, X., 2021. A crown morphology-based approach to individual tree detection in subtropical mixed broadleaf urban forests using UAV LiDAR data. *Remote Sens.* 13 (7), 1278. <http://dx.doi.org/10.3390/rs13071278>.
- Yang, M., Mou, Y., Liu, S., Meng, Y., Liu, Z., Li, P., Xiang, W., Zhou, X., Peng, C., 2022. Detecting and mapping tree crowns based on convolutional neural network and Google Earth images. *Int. J. Appl. Earth Obs. Geoinf.* 108, 102764. <http://dx.doi.org/10.1016/j.jag.2022.102764>.
- Yang, L., Wu, X., Praun, E., Ma, X., 2009. Tree detection from aerial imagery. In: *Proceedings of the 17th ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems*. pp. 131–137. <http://dx.doi.org/10.1145/1653771.1653792>.
- Zamboni, P., Junior, J.M., Silva, J.d.A., Miyoshi, G.T., Matsubara, E.T., Nogueira, K., Gonçalves, W.N., 2021. Benchmarking anchor-based and anchor-free state-of-the-art deep learning methods for individual tree detection in rgb high-resolution images. *Remote Sens.* 13 (13), 2482. <http://dx.doi.org/10.3390/rs13132482>.
- Zhang, L., Lin, H., Wang, F., 2022. Individual tree detection based on high-resolution RGB images for urban forestry applications. *IEEE Access* 10, 46589–46598. <http://dx.doi.org/10.1109/ACCESS.2022.3171585>.
- Zhu, L., Zhao, Z., Lu, C., Lin, Y., Peng, Y., Yao, T., 2019. Dual path multi-scale fusion networks with attention for crowd counting. <http://dx.doi.org/10.48550/arXiv.1902.01115>, CoRR abs/1902.01115.
- Zörner, J., Dymond, J.R., Shepherd, J.D., Wisner, S.K., Jolly, B., 2018. LiDAR-based regional inventory of tall trees—Wellington, New Zealand. *Forests* 9 (11), 702. <http://dx.doi.org/10.3390/f9110702>.