

Predicting urban tree cover from incomplete point labels and limited background information

Hui Zhang, Ankit Kariryaa, Venkanna Babu Guthula, Christian Igel, Stefan Oehmcke
 {huzh,ak,vegu,igel,stefan.oehmcke}@di.ku.dk
 University of Copenhagen
 Dept. of Computer Science
 Copenhagen, Denmark

ABSTRACT

Trees inside cities are important for the urban microclimate, contributing positively to the physical and mental health of the urban dwellers. Despite their importance, often only limited information about city trees is available. Therefore in this paper, we propose a method for mapping urban trees in high-resolution aerial imagery using limited datasets and deep learning. Deep learning has become best-practice for this task, however, existing approaches rely on large and accurately labelled training datasets, which can be difficult and expensive to obtain. However, often noisy and incomplete data may be available that can be combined and utilized to solve more difficult tasks than those datasets were intended for.

This paper studies how to combine accurate point labels of urban trees along streets with crowd-sourced annotations from an open geographic database to delineate city trees in remote sensing images, a task which is challenging even for humans. To that end, we perform semantic segmentation of very high resolution aerial imagery using a fully convolutional neural network.

The main challenge is that our segmentation maps are sparsely annotated and incomplete. Small areas around the point labels of the street trees coming from official and crowd-sourced data are marked as foreground class. Crowd-sourced annotations of streets, buildings, etc. define the background class. Since the tree data is incomplete, we introduce a masking to avoid class confusion.

Our experiments in Hamburg, Germany, showed that the system is able to produce tree cover maps, not limited to trees along streets, without providing tree delineations. We evaluated the method on manually labelled trees and show that performance drastically deteriorates if the open geographic database is not used.

KEYWORDS

urban environment, deep learning, point labels, sparse labels, crowd-sourced datasets, tree mapping, CNN, tree cover

ACM Reference Format:

Hui Zhang, Ankit Kariryaa, Venkanna Babu Guthula, Christian Igel, Stefan Oehmcke. 2023. Predicting urban tree cover from incomplete point

labels and limited background information. In *1st ACM SIGSPATIAL International Workshop on Advances in Urban-AI (UrbanAI '23)*, November 13, 2023, Hamburg, Germany. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3615900.3628791>

1 INTRODUCTION

Trees are a vital component of our ecosystems. They are vital for sustaining the biodiversity of various lifeforms and provide important services such as food, shelter, and shade [1]. In an urban setting, trees also offer benefits for physical and mental health [32]. However, trees are also highly vulnerable to change in climatic conditions [2]. Increase in global temperatures is associated with an increased global tree mortality rate, which reduces the ecosystem functioning and impacts their role in carbon storage [24]. Monitoring the amount of trees is therefore vital to devise mitigation and adaptation measures against climate change.

In this paper, we propose a method to train a deep learning model to predict tree cover in an urban setting with sparse and incomplete labels. This work differs from existing studies in several key aspects. First, our work is unique in combining several different open data sources. To the best of our knowledge, no previous study has evaluated the potential of authority-managed tree records and crowd-source annotations from an open geographic database for tree mapping. Second, we focus on urban areas, which are relatively under-explored in other work [16, 29], although many free data sources exist. Third, existing work relies on strong preprocessing and fully annotated data in which the object has either been accurately delineated [7, 20] or been annotated by a bounding box or at least a point label. An example of point labels is done by Ventura *et al.* [30], who manually annotated 100 000 trees from eight cities in the USA and collected multiple years of imagery. Also, Beery *et al.* [4] incorporated different sources of public data sets, but required multiple steps of data cleaning, resulting in nearly half of the tree records being removed. In contrast, we exclusively utilize freely available data, both for input imagery and labels, which requires no annotation efforts for training.

It is important to acknowledge that combining different sources of public data presents unique challenges, such as imbalanced classes and noisy labels, given that these data are not originally designed to be used together (see Fig.1). To make full use of the incomplete and sparsely labeled tree data as well as reduce the uncertainty of the background class, we proposed a mask regime that carefully selects pixels of trees and background with high probability of being that class. With this mask regime, we show that our approach is able to utilize this newly conjuncted dataset to predict urban trees with a balanced accuracy of 82% on sparsely labeled data

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

UrbanAI '23, November 13, 2023, Hamburg, Germany

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.
 ACM ISBN 979-8-4007-0362-1/23/11...\$15.00
<https://doi.org/10.1145/3615900.3628791>



Figure 1: Aerial image from Hamburg, Germany with street trees overlaid in red. Street trees dataset is incomplete since it does not contain any information about the trees on private land, public parks, forests, or farms.

and 84% on fully annotated data. We also introduce an objectness prior in the loss function inspired by weak supervision literature. Originally proposed in [3], pseudo-labels are derived from model predictions that are pretrained on another dataset with the same task. We derive pseudo-labels from an adapted watershed algorithm [15] to increase the extent of the object being sensed by the model for point-level supervision without requiring a pretrained model on the same task. Unlike common weak supervision scenarios that assumes sparse but fully annotated data, our incomplete annotations can lead to an incorrect objectness prior. Consequently, we also applied our mask regime to the objectness prior and restrict learning of the target class to the area close to our tree labels. This ablation study showed that the masking regime is always beneficial, while the inclusion of an objectness prior is highly dependent on its quality.

To summarize, our main contributions are as follows:

- A novel masking approach for combining noisy crowd sourced data with precise point labels.
- A dataset created from publicly available data, bringing forward the challenge of incomplete and sparse labels as well as a hand-delineated test set.
- An evaluation and comparison of different techniques to include the novel masking scheme.

2 RELATED WORK

As in many other fields, deep learning learning models have become state-of-the-art method for mapping trees and tree cover in aerial, satellite and LiDAR imagery. However, training these models in a supervised learning setting requires large volumes of manually annotated data, which is often tedious and expensive to create and requires domain expertise. Typically, these methods are trained with dense labels, such as full delineations of trees. For example, [7] manually annotated 89 899 trees on very high-resolution satellite imagery for training their deep learning models.

Recent research shows that semi- and weakly supervised learning have made great progress in the semantic segmentation of images [35]. Weakly supervised learning aims to learn from a limited amount of labels in comparison to the entire image [31, 33, 35].

Other works distinguish between different levels of weakly supervised annotations, such as bounding boxes [9], scribbles [17], points [15, 34], image labels [23], pixel-level pseudo labels generated with class activation maps [12, 25, 27], and also a text-driven semantic segmentation [18]. While fully-labelled data is limited, point labels are also used in instance segmentation methods, such as [13] introduced a novel learning scheme in instance segmentation with point labels and [14] proposed point-level instance segmentation with two branch network such as localisation and embedding branch.

Interactive segmentation with point labels started a few decades back and is still an active research topic [6]. These segmentation models started training with point labels that annotate entire objects. Lin *et al.* [17] proposed ScribbleSup based on a graphical model that jointly propagates information from scribble and points to unmarked pixels and learns network parameters without a well-defined shape. Maninis *et al.* [19] proposed a framework with a point-level annotation that follows specific labeling instructions such as left-most, right-most, top, and bottom pixels of segments. Bearman *et al.* [3] proposed a methodology by incorporating objectness potential in the training loss function in segmentation models with image and point-level annotations. Li *et al.* [15] utilised an objectness prior similar to [3] but instead of a convolutional neural network (CNN) output they utilize distances in the pixel and colour space, meaning that the further away in the image and the more different the colour, the objectness decreases. Zhang *et al.* [34] proposed a contrast-based variational model [22] for semantic segmentation that supports reliable complementary supervision to train a model for histopathology images. Their proposed method incorporates the correlation between each location in the image and annotations of in-target and out-of-target classes. The weak supervision part of our research is inspired by [3, 15, 34], as we have only a single point for each tree, we use point labels in combination with denser background information while considering an objectness prior.

In contrast to these scenarios, we consider the added challenge of incomplete annotations, meaning some relevant objects in an image might not be annotated at all.

3 OPENCITYTREES DATASET

Public agencies often maintain valuable records of trees and other public attributes such as roads, parking areas, buildings, etc. These datasets are, however, often noisy due to differences in collection techniques, lack of the common data collection standards, noisy sensors, and lack of records of temporal changes. Moreover, they are mostly not developed for the goal of training supervised deep learning models or for use in conjunction with other modalities such as aerial or satellite imagery. As such, they are potentially underutilized in research. To demonstrate their usefulness, we created a new dataset for weakly supervised segmentation from such records.¹

3.1 Input images

To demonstrate the usefulness of public but incomplete datasets, we use the aerial images from Hamburg, Germany as input for our

¹<https://doi.org/10.17894/ucph.b1aa4ca2-9a4b-40d0-aa87-c760e69bf703>

models. These images contain 3 channels (RGB) at a 0.2 m/pixel resolution and they were downloaded from the data portal of the Spatial Data Infrastructure Germany (SDI Germany)². The images were captured in May 2016. As seen in Fig. 1, the individual features such as trees, buildings, and cars on the streets are visible to human eyes. We downloaded 27 image tiles of 5000 × 5000 pixels (*i.e.* covering a 1 km × 1 km area) within the bounds (9.9479070E, 53.4161744N, 9.9684731E, 53.6589539N). These tiles extend from the north to the south border of Hamburg but are limited to 1 km strip close to the city center. Hamburg is situated on the coast of the Elbe river with a densely populated city-center. Along its border (*i.e.* away from the city center), the city-state also contains suburbs, farms and forested area. The chosen images capture all these different characteristics of the city along the north-south gradient. Since images are captured in early spring, many of the trees are without leaves, making certain trees more challenging to identify.

When designing the dataset, we considered that additional height data derived from LiDAR could potentially enhance results. However, we decided against including it because there are practical constraints associated with LiDAR data availability and collection. High-resolution LiDAR data (*e.g.*, submeter similar to our RGB source) often remains inaccessible due to regulatory limitations, especially concerning drone or plane flights over urban areas. Additionally, acquiring LiDAR measurements is more costly compared to RGB measurements, which could make frequent temporal analyses of urban tree cover infeasible. For instance, the open data portal we used, does not have submeter height measurements available for Hamburg.

3.2 Label data

Two sources of labeled data are combined:

Ground truth for trees. The Authority for Environment and Energy of the city of Hamburg maintains a list of all street trees³ as recorded on the 6th of January 2017. The dataset contains various attributes of individual trees such as location, height, width, species, age, and condition. However, as the name suggests, this information is limited to the trees along the streets of Hamburg and does not include information about trees on private land, in public parks, or in forested areas. Unlike other data usually used in point-supervision where each object is assumed to be annotated with at least one point, we have incomplete annotations, increasing the ambiguity of the background class. In the area of interest, the dataset contains information about 11 366 trees. These trees are from 136 unique species. In Fig. 1, trees in the street trees dataset are overlaid in red circles. Each tree is provided as a point referenced in a local reference system (EPSG:32632 - WGS 84). However, the point location of the tree label can be inaccurate, for example, it might not overlap with the center of the tree or, in the worst case, any part of the trees due to the geo-location errors. Another challenge with the dataset is that distribution of species of the street trees may vary significantly from the distribution of trees species in forests, parks, farms, or gardens.

As a second source of ground truth data, we use OpenStreetMap (OSM) [10]. Within these bounds and the tag 'natural': 'tree', OSM

contains the location of 6375 trees. Out of these, 145 trees contain species information (24 unique species). These OSM data offer information regarding trees in private and public areas along with street trees.

Table 1: Description of objects defining the non-tree class. The buffer distance is in meters and negative buffers shrinks the object. Only vectors with non-negative area were chosen.

Type	OSM tag	Count	Buffer	Comments
Buildings	'building':True	23 075	-5m	Buildings of all types
Roads	'highway':True	111	-7m	Mostly around parking areas or bus terminals
Sports	'leisure': 'pitch'	135	-7m	Soccer pitches and similar types of grass surfaces

Ground truth for non-trees classes. Table 1 provides an overview of the objects that we use to define the non-tree class. The non-tree classes are mostly dominated by buildings which provide relevant information about different construction material and roof types. While the area contains abundant roads, it is a tricky class to consider for the true negatives since the trees are often planted next to the roads and large parts of tree canopies overlap with roads. We only used road data if they had an associated area (*i.e.* stored as polygon or multi-polygon). We used OSMNX library to download data from OSM [5]. Sports pitches, which includes grassed surfaces such as soccer pitches, are limited to 135 instances and it is only classes that provides information on grass which is easy to confuse with trees.

3.3 Challenging aspects of the data

By combining tree inventories and geographic data from existing public records, we create a rich dataset, without the need for additional acquisition of labor-intensive annotations. Public records maintain valuable information about trees and other public attributes. However, using incomplete public records for tree prediction also introduces a number of challenges:

Sparse labels: The ground truth of trees are given as point labels that cover most public streets, some public parks, and a few private places. These annotations are incomplete and only represent a small portion of urban trees. In addition, these street trees are also sparsely distributed.

Presence of noise: Although the tree census data and aerial images are obtained from relatively close point of times, it is important to note that changes in the tree population might have occurred during the time gap. Trees could have been removed, died, or new trees might have been planted. Besides, there are geo-location errors as mentioned before, any nearby pixel of the tree could be labeled as the tree centroid.

Image quality: The quality of aerial imagery can vary for different tree species. The images were captured in early spring, when deciduous trees have not yet grown leaves. In addition, renewal of growth in trees near streets may be influenced by extended period

²<https://www.geoportal.de/portal/main/>

³<http://suche.transparenz.hamburg.de/dataset/strassenbaumkataster-hamburg7>

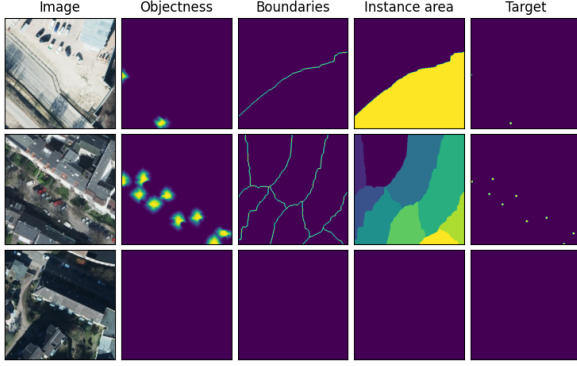


Figure 2: Objectness prior maps (col 2) and instance areas (col 4) were generated using input images (col 1) and locations of tree centroids (col 5) according to [15]. The boundaries (col 3) are derived where instance areas are touching.

of illumination and emissions from the streets [21]. As a result, these trees may not be well represented in the aerial image.

Invisible trees: There are trees located within shadows of nearby tall buildings, darkening the image and increasing potential class confusion between tree and shadows.

4 LEARNING FROM INCOMPLETE & SPARSE LABELS

Our main challenge in training a tree segmentation model is obtaining accurate and effective labels. Coming from open-data sources, however, labels are incomplete, meaning that not all trees or non-tree objects in an area are annotated. These incomplete labeling deviates from the typical definition of weakly supervised learning [3, 15], where we assume sparse labels (e.g., points, scribbles, bounding boxes, ...) are available for every relevant object. In addition, tree labels are only available on a point level, meaning a single point represents a tree although the tree canopy encompasses a larger area. The non-tree labels are taken from OSM thus describing only parts of the image, in addition, we chose to shrink their shape to avoid overlap with potential tree that are not covered by our dataset (e.g., a tree reaching over a building), see Table 1.

We frame the learning task as binary semantic segmentation of trees and introduce concepts to deal with the incomplete sparse labels. To that end, we consider a training set

$$T = \{(x_1, y_2), \dots, (x_n, y_n)\} \subset X \times Y$$

with images $x_i \in X = \mathbb{R}^{w \times h \times c}$, a segmentation mask $y_i \in Y = \{-1, 1\}^{w \times h}$, and number of samples n . Further, w , h , and c correspond to the width, height, and number of input channels, respectively. The pixels containing the non-trees objects are considered as the negative class samples. In our training dataset, we treat the pixels in a 60 cm radius (7×7 pixels) around the point coordinate of a tree as positive class labels, which increases the number of positive training labels substantially.

Training in such a setting is non-trivial. For example, learning a semantic segmentation only given point labels is challenging

because information about the spatial expand of the objects in question is limited. Previous research introduces this spatial expand information by means of an objectness prior. The objectness prior gives an estimate of the class likelihood per pixel. As shown in Figure 2, given the location of the trees, the algorithm estimates the potential spatial extent for each tree.

This prior can come from pretrained models on similar tasks [3], but also from classic algorithms, e.g. inspired by watershed segmentation [15]. Our approach uses these two loss functions in conjunction:

$$\mathcal{L} = \mathcal{L}^{\text{sup}}(f(x) \odot m, y \odot m) + \mathcal{L}^{\text{obj}}(f(x) \odot r, o \odot r, \hat{r} \odot r) \cdot \beta, \quad (1)$$

where \mathcal{L}^{sup} is the supervised loss (e.g., binary cross entropy (BCE) loss) that learns from the labeled data and the objectness loss \mathcal{L}^{obj} , where this prior information is utilized. Here, \odot denotes the selection operator that chooses elements where the learning mask m is set to 1 and returns the elements as a flattened vector. The parameters of \mathcal{L}^{obj} are the predictions $f(x)$, the objectness prior o , and an instance region \hat{r} . Note, since no pre-trained CNN [3] on tree segmentation was easily available, we utilized the method described by [15] to calculate o and \hat{r} . To obtain o , we calculate the distance matrix $\Delta \in \mathbb{R}^{w,h}$ by applying the adjusted watershed algorithm [28] as in [15] with the point labels being used as markers and then transforming these distances into a pseudo-probability distribution $o = e^{-\alpha \Delta^2}$, with $\alpha = 10$ to create fast decay of values the farther away from an actual label. The current settings for these pseudo-probabilities were explored during a preliminary study on the training set but the ones provided by [15] turned out to perform best. From the same adjusted watershed output, we use the watershed instance assignments as \hat{r} . β is trade-off parameters to change the influence of the objectness loss. See Figure 2 for an exemplary input and objectness-related attribute. In our incomplete label setting, the generated objectness can only capture the trees indicated by point labels. Therefore, to represent where labels are available, we declare two learning masks $m \in \{0, 1\}^{w \times h}$ and $r \in \{0, 1\}^{w \times h}$, where 1 means a label is present and 0 corresponds to missing label information. These masks can be defined in several ways as we explore in the experimental section and can be considered one of main contributions of this paper.

For the objectness loss, we extend the binary cross-entropy similarly to [3, 15]

$$\mathcal{L}^{\text{obj}} = -\frac{1}{|o|} \sum_{i=1}^{|o|} \text{BCE}(f(x) \odot \hat{r}_i, o \odot \hat{r}_i), \quad (2)$$

where each tree instance is calculating its own loss value depending on the instance region $\hat{r} \in \{0, 1\}^{w \times h}$. The number of tree instance $|\hat{r}|$ changes for each sample, as does the number of pixels in each region. Averaging inside the instance sum effectively weights each instance the same, regardless of size.

5 EXPERIMENTAL EVALUATION IN HAMBURG

5.1 Ablation study

The choice of masks m and r is crucial in our sparse and incomplete label setting. To that end, we compared five different training scenarios as shown in Table 2. The *baseline* scenario is only using the supervised loss without any masking. The public authority and OSM tree labels are expanded from a point to a disk m^{disk} of radius 1.5 m, which is indicated by $m^{\text{disk}} = 1$. The second scenario, called *Obj*, uses the objectness loss over the entire image in combination with supervised loss and we mask out all pixels with negative labels except on the boundaries of the instance region \hat{r} (see Figure 2). *Obj* is a reimplementation of [15]. The third scenario uses the supervised loss along with our proposed masking scheme, termed *Mask*. Here we do not consider the objectness loss and only evaluate the supervised loss where we have positive labels (indicated by $y = 1$), and where we have information about the shrunken OSM non-tree objects m^{OSM} , which is indicated by $m^{\text{OSM}} = 1$. In addition, in shrunken OSM non-tree objects, we remove negative pixels that are within 1.5 m of a positive label. In the fourth scenario, we combine our masking approach with objectness in *MaskObj*, by employing the objectness loss but restricting it to the 1.5 m radius around the positive labels. Lastly, we add an additional constraint to the objectness by ignoring all the pseudo-probabilities that are below 0.2, which will reduce the learning about the negative class in \mathcal{L}^{obj} , which we refer to as *MaskObjThresh*.

5.2 Network architecture, loss function, and hyperparameters

To address the tree segmentation task, we employ a fully-convolutional network based on the U-Net architecture [26].

Experimental settings. Among other things, in the past U-Nets have been used for semantic segmentation of trees in satellite imagery [7]. We adapted the U-Net architecture by applying batch normalization [11] instead of dropout layers and replacing ReLU with ELU [8] as activation functions. We use binary cross entropy (BCE) as loss function as our supervised error measure. The aerial images were split into 300×300 patches and a batch consists of 36 patches. For training and hyperparameter optimization, the dataset was split into 80% training set (3566 patches), 20% validation set (788 patches). To improve training stability, we accumulated gradients over 14 batches (i.e 504 images) before the optimizer step. The model is trained for 500 epochs and the final weights were chosen w.r.t. the best recall score on the validation set.

Evaluation on sparse and dense labels. The evaluation of the model’s performance was done with two types of data annotations, point annotation, and dense object annotation. These two datasets are spatially independent. First, we evaluated the model on the point-annotated data from 28 tiles (4169 patches) within the bounds (9.962748E, 53.407065N, 9.83603E, 53.658832N), which is a 1 km \times 28 km stripe adjacent to the training data stripe. None of these tiles were used for training or intra-model validation and the ground truth dataset for them was created in exactly the same way as described in Section 3. None of the pseudo labeling (e.g., extending of point labels to 4 pixels or a disk as label) was utilized

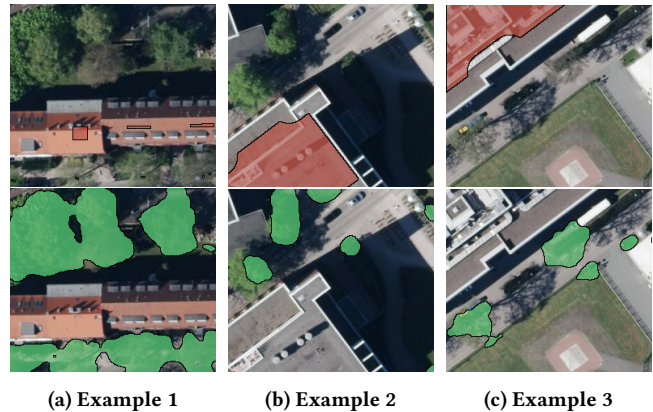


Figure 3: Three examples of target (top) and possible predicted segmentation (bottom). The predicted positive class is overlaid in green. In the target examples, a red overlay indicates the negative class and transparency means the learning mask is 0. For the predicted segmentation, only the positive class is shown and the negative class is transparent.

during evaluation, meaning that for point labels only the corresponding pixel is considered and for the background class only the negatively buffered area. The sparse street tree dataset and OSM had information on 14 137 trees within these bounds.

To evaluate our models performance on dense object prediction, we manually annotated a tile within a 1 km² area, which is 3 km to the east of our training data. The delineation work was done using QGIS and is mainly based on the input image, which was cross-referenced with Bing and Google Satellite Maps. The annotation was then verified within the authors’ group, which eliminates some bias. To utilize this dataset for an unbiased tree cover estimate, we split it further into a model selection set and a test set. We applied only the best model from the model selection set in terms of IoU to the test set.

5.3 Sparse Label Results

The results of the sparse label test set are given in Table 2. It is crucial to acknowledge the highly imbalanced nature of the dataset when evaluating with sparse labels. Due to this significant imbalance, the number of false positives can be far greater than the true positive, leading to a substantially low precision value. Specifically, due to the class imbalance with 2448 times more negative class pixels than positive class pixels, the precision of our models was only around 3%. Therefore we focus on the recall (sensitivity) of the target class and balanced accuracy (BA) to evaluate the model performance on sparse labels.

The *baseline* model performed worst and appears to mainly predict the background class. Performing best was the *Mask* model w.r.t. recall with 90% and *MaskObj* w.r.t. BA with 84%. Even though the BA of *Obj* is close to the mask models with 78%, the recall value is comparatively low with 59%. In Figure 3 exemplary target and prediction segmentation masks are shown.

Table 2: Results on sparse and delineated data across different ablation settings. For model selection set of the delineated labels, we compare intersection over union of the tree (IoU) of the tree class, F1, and balanced accuracy (BA) scores. The sparse labels are compared w.r.t. their recall and BA scores. Additional masks are the shrunken OSM non-tree objects $m^{\text{OSM}} \in \{0, 1\}^{w \times h}$, a 1.5 m disk $m^{\text{disk}} \in \{0, 1\}^{w \times h}$ around each positive values in y , the bounds b between instances derived from the instance region map r , and a mask of ones $1^{w \times h} \in \{1\}^{w \times h}$. Results of the *baseline* model were not calculated for the delineated data set since the model was discarded due to the sparse label performance.

Name	Description	Setting	Delineated			Sparse	
			IoU _d	F1 _d	BA _d	Recall _s	BA _s
<i>baseline</i>	Neither masking nor objectness and enlarging positive labels to a circle with 1.5 m radius.	$\beta = 0$ $m = 1^{w \times h}$ $y = m^{\text{disk}}$	—	—	—	0.0005	0.5002
<i>Obj</i>	Reimplementation of [15].	$\beta = 1$ $m = y \cup b$ $\hat{r} = 1^{w \times h}$	0.1191	0.5479	0.5575	0.5908	0.7844
<i>Mask (ours)</i>	Only supervised loss and masking out unknown areas.	$\beta = 0$ $m = y \cup (m^{\text{OSM}} \setminus m^{\text{disk}})$	0.4839	0.7551	0.8119	0.8994	0.8205
<i>MaskObj (ours)</i>	As [15] but restricting objectness loss to 1.5 m radius around points.	$\beta = 1$ $\hat{r} = m^{\text{disk}}$ $m = y \cup (m^{\text{OSM}} \setminus m^{\text{disk}})$	0.3364	0.7289	0.6978	0.7771	0.8399
<i>MaskObjThresh (ours)</i>	As MaskObj but removing objectness smaller than the specified threshold ($t = 0.2$).	$\beta = 1$ $\hat{r} = m^{\text{disk}} \cap (o \geq t)$ $m = y \cup (m^{\text{OSM}} \setminus m^{\text{disk}})$	0.4805	0.7660	0.7870	0.8345	0.8135

5.4 Delineation Results

The intersection over union (IoU), F1, and BA score on the model selection set can be seen in Table 2. We omitted the *baseline* because of the subpar results on the sparse test set. The class imbalance changes since these annotation are fully delineated, particularly, we now have 3.42 negative pixels for one positive pixel. This change makes the use precision viable, which is why we consider the F1-score. The original masking model *Mask* performed best in IoU and BA, even though the difference in IoU compared to *MaskObjThresh* is marginal. *Obj* only shows a BA of 56% which is much smaller than the 78% on the sparse data, showing a lack of generalization for this approach. Using any kind of masking scheme seems to improve the results.

In Figure 4 we show the normalized and unnormalized confusion matrix results on the model selection set. Note that *Obj* and *MaskObj* underpredicts the positive class, which corresponds to the low accuracy and a false-negative rate. For *MaskObjThresh* and *Mask* the accuracy of the positive class is considerably higher, even though the false positive rate increases. Interestingly, thresholding the objectness prior improves performance compared to the model without, indicating that the prior hold misleading information regarding the background class.

Based on the results on the model selection set, we decided that the best model is the one without any objectness loss but with masking (*Mask*). This decision is based on the better metrics but also on the simpler learning setup. Our masking regime only requires a sensible choice for the excluded areas, while the objectness prior

additionally requires a choice of how to calculate it. For delineated test set the *Mask* model achieved an IoU of 0.4253, F1 score of 0.7393 and BA of 83.63%. These results are comparable to the other sets, indicating a good generalization performance.

Finally, the predictions on the dense label set are shown in Figure 5. Within the bounds of the test area (591 609 m²), we detected a total tree cover of 177 142 m² (29.9%) compared to the annotated 96 946 m² (16.4%). This shows an overestimation but as seen in the map, some of the trees were missed in the annotated dataset or were difficult to delineate without ambiguity.

5.5 Discussion of Ablation results

Baseline. The baseline segmentation model trained on incompletely labeled data with ambiguous information about the target class and the background class is not able to learn tree features thus predicts all pixels as background. As shown in the first row of Table 2, the recall of the target class is close to 0 when evaluated on the sparse labels. The model mainly predicts the background class, which aligns with our assumption that the class imbalance without our masking is challenging to overcome.

Effect of Objectness Prior. By introducing the objectness prior in the *Obj* model, pixels spatially and chromatically close to the tree centroid are given higher probability of being that tree. Effectively, it expands the learning from only known tree pixels to also include possible pixels. However, our task is different from [15] because many objects are unlabeled, which makes the boundaries generated from instance areas less representative of the background class.

True Label	Non-tree	0.79 38.5e+5	0.21 10.4e+5	0.99 48.6e+5	0.01 0.3e+5
	Tree	0.16 2.3e+5	0.84 12.0e+5	0.88 12.6e+5	0.12 1.7e+5
		(a) Mask		(b) Obj	
True Label	Non-tree	0.96 47.1e+5	0.04 1.8e+5	0.86 41.9e+5	0.14 7.1e+5
	Tree	0.62 8.8e+5	0.38 5.4e+5	0.28 4.0e+5	0.72 10.3e+5
		(c) MaskObj		(d) MaskObjThresh	
		Non-tree	Tree	Non-tree	Tree
		Predicted		Predicted	

Figure 4: Confusion matrix across different ablation settings, created from the delineated model selection dataset. The normalized confusion matrix is shown in bold font (top number) and through color, while the lower number represents absolute number of samples. (a) Mask: Our initial method that utilizes masking, (b) Obj: Model as presented in [15], (c) MaskObj: Combining masking and objectness, (d) MaskObjThresh: Restricting learning to positive labels for objectness loss. We do not show the baseline model here, because the model mainly predicted the negative class (e.g., in the first column both values are close to 1).

Our results show that objectness helps when training in a weakly supervised setting with imbalanced data, but the incomplete labeling is not accounted for in this case and the models including our masking regime (*MaskObj* and *MaskObjThresh*) outperform the *Obj* model.

Effect of Mask. The masking regime is crucial in an imbalanced and sparsely labeled data setting. Without masking, the baseline model failed to predict the target class completely. As shown in Table 2, the mask regime introduced more performance improvements than the objectness prior, with a significantly improved IoU score for models with masking. Among those models, the *Mask* model achieved the highest IoU and BA on the delineated data and best recall score on the sparse data. In summary, the mask regime is comparatively simple to implement, costs less computational resources, and requires less fine tuning than using the objectness prior.

Effect of combining Objectness Prior and Mask. The mask regime consistently improved the performance of the models learning from the objectness prior, but there was no improvement compared to the model applying masking without the objectness prior.

A possible reason for this might be that in our complicated urban environment setting, the spatial context of trees might vary quite a bit (e.g., small and large trees) and there is a chance that the objectness map would highlight the object around a tree that is actually not a tree (see Figure 2). Since the extend and cutoff of the probabilities depend on hyperparameters when creating the prior, there could be settings giving a good performance for one but not all the different scenarios. This brings us to the conclusion that the objectness prior in its current form does not yield any benefits compared to simply applying the mask regime on its own.

6 CONCLUSION

In this paper, we create a new tree segmentation dataset from public data to train a deep learning model for semantic segmentation of urban trees. This dataset is challenging because it consists of incomplete and sparse point labels for trees and carefully selected background objects from OSM. We address this label incompleteness and sparsity by proposing a loss masking regime into our model design including domain knowledge. Further, we expand on the weakly supervised technique of learning from objectness priors by utilizing the same masking regime.

Our evaluation shows that the best performance was achieved when only using our masking regime, with a test performance on a fully delineated set of 0.43 IoU and 84% balanced accuracy. This indicates that the chosen objectness prior was not helpful for this task while our mask regime is beneficial when dealing with incomplete and sparsely annotated data. This even holds when combining it with other approaches, such as the *Obj* model. Besides, our mask regime is simple to implement, lower in computational resource requirements, and requires less fine-tuning of hyperparameters. While including *Obj* or its variants requires to identify and compute an appropriate objectness prior o for each new task and training sample. This overhead is an inherent aspect of the objectness methods. Our results on both point labels and a manually delineated evaluation set demonstrates the hidden potential of public datasets for mapping urban trees.

In the future, we will investigate the usefulness of self-supervised networks pretrained on large datasets in conjunction with weak labels to further improve the mapping performance. Evaluation in other urban areas would also be of interest to validate the generalization performance. Since our current masking prior is calibrated to our Hamburg dataset, it may not transfer well to new areas. Therefore we would like to investigate if the masking prior could be learned or adjusted from unsupervised or semi-supervised networks.

ACKNOWLEDGMENTS

This work was supported by the research grant DeReEco from VILLUM FONDEN (grant number 34306), the PerformLCA project (UCPH Strategic plan 2023 Data+ Pool), and the grant “Risk-assessment of Vector-borne Diseases Based on Deep Learning and Remote Sensing” (grant number NNF21OC0069116) by the Novo Nordisk Foundation.

REFERENCES

- [1] Raf Aerts and Olivier Honnay. 2011. Forest restoration, biodiversity and ecosystem functioning. *BMC Ecology* 11, 1 (2011), 1–10. 1



Figure 5: Results on the fully delineated dataset. Red polygons are the manually annotated ground truth. Validation data is shown within the red bounds, which is used for the model selection. As discussed in Section 5, the *Mask* model is the final model chosen for prediction, which shows as the yellow mask on top of the satellite image. The segmentation mask in green are from the model *MaskObjThresh*. The area bounded by the green lines shows the test set. Within the test area, we only apply *Mask*. In cutout A the model predicts one of the trees correctly and interestingly does not predict the shadow of said tree. This tree is also of bright color since it is in bloom making it a hard example. However, the second tree in top left corner was not found, possibly because the coloration is close to the one of shadows. After manual rechecking of the predicted tree in left lower corner, we found that there is a tree not annotated originally (e.g., a false negative label). Cutout B is interesting because it contains trees without leaves. The model annotated these trees effectively and even finds another tree that was missed by the annotators (bottom right). This also demonstrates how difficult it is to define clear boundaries for these trees.

- [2] Jean-Francois Bastin, Yelena Finegold, Claude Garcia, Danilo Mollicone, Marcelo Rezende, Devin Routh, Constantin M Zohner, and Thomas W Crowther. 2019. The global tree restoration potential. *Science* 365, 6448 (2019), 76–79. 1
- [3] Amy L. Bearman, Olga Russakovsky, Vittorio Ferrari, and Li Fei-Fei. 2016. What's the Point: Semantic Segmentation with Point Supervision. In *European Conference on Computer Vision (ECCV) (Lecture Notes in Computer Science, Vol. 9911)*. Springer, 549–565. https://doi.org/10.1007/978-3-319-46478-7_34 2, 4
- [4] Sara Beery, Guanhang Wu, Trevor Edwards, Filip Pavetic, Bo Majewski, Shreyasee Mukherjee, Stanley Chan, John Morgan, Vivek Rathod, and Jonathan Huang. 2022. The Auto Arborist Dataset: A Large-Scale Benchmark for Multiview Urban Forest Monitoring Under Domain Shift. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 21294–21307. 1
- [5] Geoff Boeing. 2017. OSMnx: New methods for acquiring, constructing, analyzing, and visualizing complex street networks. *Computers, Environment and Urban Systems* 65 (2017), 126–139. 3
- [6] Yuri Y. Boykov and Marie-Pierre Jolly. 2001. Interactive graph cuts for optimal boundary & region segmentation of objects in N-D images. In *International Conference on Computer Vision (ICCV)*, Vol. 1. 105–112 vol.1. <https://doi.org/10.1109/ICCV.2001.937505> 2
- [7] Martin Brandt, Compton Tucker, Ankit Kariryaa, Kjeld Rasmussen, Christin Abel, Jennifer Small, Jerome Chave, Laura Vang Rasmussen, Pierre Hiernaux, Abdoul Aziz Diouf, Laurent Kergoat, Ole Mertz, Christian Igel, Fabian Gieseke, Johannes Schöning, Sizhuo Li, Katherine Melocik, Jesse Meyer, Scott Sinno, Eric Romero, Erin Glennie, Amandine Montagu, Morgane Dendoncker, and Rasmus Fensholt. 2020. An unexpectedly large count of trees in the West African Sahara and Sahel. *Nature* 587, 7832 (2020), 78–82. 1, 2, 5
- [8] Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. 2016. Fast and accurate deep network learning by exponential linear units (ELUs). 5
- [9] Jifeng Dai, Kaiying He, and Jian Sun. 2015. BoxSup: Exploiting Bounding Boxes to Supervise Convolutional Networks for Semantic Segmentation. In *International Conference on Computer Vision (ICCV)*. IEEE, 1635–1643. <https://doi.org/10.1109/ICCV.2015.191> 2
- [10] Mordechai Haklay and Patrick Weber. 2008. Openstreetmap: User-generated street maps. *IEEE Pervasive Computing* 7, 4 (2008), 12–18. 3
- [11] Sergey Ioffe and Christian Szegedy. 2015. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning (ICML)*. PMLR, 448–456. 5
- [12] Dahyun Kang, Piotr Koniusz, Minsu Cho, and Naila Murray. 2023. Distilling Self-Supervised Vision Transformers for Weakly-Supervised Few-Shot Classification & Segmentation. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 19627–19638. 2
- [13] Beomyoung Kim, Joonhyun Jeong, Dongyoon Han, and Sung Ju Hwang. 2023. The Devil is in the Points: Weakly Semi-Supervised Instance Segmentation via Point-Guided Mask Representation. [arXiv:2303.15062 \[cs.CV\]](https://arxiv.org/abs/2303.15062) 2
- [14] Issam H. Laradji, Negar Rostamzadeh, Pedro O. Pinheiro, David Vázquez, and Mark Schmidt. 2020. Proposal-Based Instance Segmentation With Point Supervision. In *International Conference on Image Processing (ICIP)*. IEEE, 2126–2130. <https://doi.org/10.1109/ICIP40778.2020.9190782> 2
- [15] Shijie Li, Neel Dey, Katharina Bermond, Leon von der Emde, Christine A. Curcio, Thomas Ach, and Guido Gerig. 2021. Point-Supervised Segmentation Of Microscopy Images And Volumes Via Objectness Regularization. In *International Symposium on Biomedical Imaging (ISBI)*. IEEE, 1558–1562. <https://doi.org/10.1109/ISBI48211.2021.9433963> 2, 4, 5, 6, 7
- [16] Weijia Li, Haohuan Fu, Le Yu, and Arthur Cracknell. 2016. Deep learning based oil palm tree detection and counting for high-resolution remote sensing images. *Remote Sensing* 9, 1 (2016), 22. 1
- [17] Di Lin, Jifeng Dai, Jiaya Jia, Kaiying He, and Jian Sun. 2016. ScribbleSup: Scribble-Supervised Convolutional Networks for Semantic Segmentation. In *Computer Vision and Pattern Recognition (CVPR)*. IEEE Computer Society, 3159–3167. 2
- [18] Yuqi Lin, Minghao Chen, Wenxiao Wang, Boxi Wu, Ke Li, Binbin Lin, Haifeng Liu, and Xiaofei He. 2023. CLIP Is Also an Efficient Segmenter: A Text-Driven Approach for Weakly Supervised Semantic Segmentation. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 15305–15314. 2
- [19] Kevis-Kokitsi Maninis, Sergi Caelles, Jordi Pont-Tuset, and Luc Van Gool. 2018. Deep Extreme Cut: From Extreme Points to Object Segmentation. In *Computer Vision and Pattern Recognition (CVPR)*. Computer Vision Foundation / IEEE Computer Society, 616–625. <https://doi.org/10.1109/CVPR.2018.00071> 2
- [20] José Augusto Correa Martins, Keiller Nogueira, Lucas Prado Osco, Felipe David Georges Gomes, Danielle Elis Garcia Furuya, Wesley Nunes Gonçalves, Diego André Sant'Ana, Ana Paula Marques Ramos, Veraldo Liesenberg, Jefferson Alex dos Santos, Paulo Tarso Sanches de Oliveira, and José Marcato Junior. 2021. Semantic segmentation of tree-canopy in urban environment with pixel-wise deep learning. *Remote Sensing* 13, 16 (2021), 3054. 1
- [21] Edwin B Matzke. 1936. The effect of street lights in delaying leaf-fall in certain trees. *American Journal of Botany* (1936), 446–452. 4
- [22] David Mumford and Jayant Shah. 1989. Optimal approximations by piecewise smooth functions and associated variational problems. *Communications on Pure and Applied Mathematics* 42 (1989), 577–685. 2
- [23] George Papandreou, Liang-Chieh Chen, Kevin Murphy, and Alan L. Yuille. 2015. Weakly- and Semi-Supervised Learning of a DCNN for Semantic Image Segmentation. [arXiv:1502.02734 \[cs.CV\]](https://arxiv.org/abs/1502.02734) 2
- [24] Hans-Otto Pörtner, Debra C Roberts, H Adams, C Adler, P Aldunce, E Ali, R Ara Begum, R Betts, R Bezner Kerr, R Biesbroek, et al. 2022. Climate change 2022: Impacts, adaptation and vulnerability. *IPCC Sixth Assessment Report* (2022). 1
- [25] Shenghai Rong, Bohai Tu, Zilei Wang, and Junjie Li. 2023. Boundary-Enhanced Co-Training for Weakly Supervised Semantic Segmentation. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 19574–19584. 2
- [26] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. 2015. U-Net: Convolutional Networks for Biomedical Image Segmentation. In *Medical Image Computing and Computer-Assisted Intervention (MICCAI) (Lecture Notes in Computer Science, Vol. 9351)*, Nassir Navab, Joachim Hornegger, William M. Wells III, and Alejandro F. Frangi (Eds.). Springer, 234–241. https://doi.org/10.1007/978-3-319-24574-4_28 5
- [27] Lixiang Ru, Heliang Zheng, Yibing Zhan, and Bo Du. 2023. Token Contrast for Weakly-Supervised Semantic Segmentation. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. 3093–3102. 2
- [28] Pierre Soille. 1999. *Morphological image analysis: principles and applications*. Vol. 2. Springer. 4
- [29] Compton Tucker, Martin Brandt, Pierre Hiernaux, Ankit Kariryaa*, Kjeld Rasmussen, Jennifer Small, Christian Igel, Florian Reiner, Katherine Melocik, Jesse Meyer, Scott Sinno, Eric Romero, Erin Glennie, Yasmin Fitts, August Morin, Jorge Pinzon, Devin McClain, Paul Morin, Claire Porter, Shane Loeffle, Laurent Kergoat, Bil-Assanou Issoufou, Patrice Savadogo, Jean-Pierre Wigneron, Benjamin Poulter, Philippe Ciaï, Robert Kaufmann, Ranga Myneni, Sassan Saatchi, and Rasmus Fensholt. 2023. Sub-continental-scale carbon stocks of individual trees in African drylands. *Nature* 615, 7950 (2023), 80–86. 1
- [30] Jonathan Ventura, Milo Honsberger, Cameron Gonsalves, Julian Rice, Camille Pawlak, Natalie LR Love, Skyler Han, Viet Nguyen, Keilana Sugano, Jacqueline Doremus, G. Andrew Fricker, Jenn Yost, and Matt Ritter. 2022. Individual tree detection in large-scale urban environments using high-resolution multispectral imagery. [arXiv preprint arXiv:2208.10607](https://arxiv.org/abs/2208.10607) (2022). 1
- [31] Alexander Vezhnevets and Joachim M. Buhmann. 2010. Towards weakly supervised semantic segmentation by means of multiple instance and multitask learning. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 3249–3256. <https://doi.org/10.1109/CVPR.2010.5540060> 2
- [32] Kathleen L Wolf, Sharon T Lam, Jennifer K McKeen, Gregory RA Richardson, Matilda van den Bosch, and Adrina C Bardekjian. 2020. Urban trees and human health: A scoping review. *International Journal of Environmental Research and Public Health* 17, 12 (2020), 4371. 1
- [33] Hongshan Yu, Zhengeng Yang, Lei Tan, Yaonan Wang, Wei Sun, Mingui Sun, and Yandong Tang. 2018. Methods and datasets on semantic segmentation: A review. *Neurocomputing* 304 (2018), 82–103. <https://doi.org/10.1016/j.neucom.2018.03.037> 2
- [34] Hongrun Zhang, Liam Burrows, Yanda Meng, Declan Sculthorpe, Abhik Mukherjee, Sarah E. Coupland, Ke Chen, and Yalin Zheng. 2023. Weakly Supervised Segmentation With Point Annotations for Histopathology Images via Contrast-Based Variational Model. In *Computer Vision and Pattern Recognition (CVPR)*. IEEE, 15630–15640. 2
- [35] Man Zhang, Yong Zhou, Jiaqi Zhao, Yiyun Man, Bing Liu, and Rui Yao. 2020. A survey of semi- and weakly supervised semantic segmentation of images. *Artificial Intelligence Review* 53, 6 (2020), 4259–4288. <https://doi.org/10.1007/s10462-019-09792-7> 2