



From simple labels to semantic image segmentation: Leveraging citizen science plant photographs for tree species mapping in drone imagery

Salim Soltani^{1,2,3*}, Olga Ferlian^{3,5}, Nico Eisenhauer^{3,5}, Hannes Feilhauer^{1,2,3,4}, and Teja Kattenborn^{1,3}

¹Remote Sensing Centre for Earth System Research (RSC4Earth), Leipzig University, Germany

²Center for scalable data analytics and artificial intelligence (ScaDS.AI), Leipzig University, Germany

³German Centre for Integrative Biodiversity Research (iDiv), Halle-Jena-Leipzig, Germany

⁴Helmholtz Centre for Environmental Research, Leipzig, Germany

⁵Institute of Biology, Leipzig University, Germany

* Corresponding author: salim.soltani@uni-leipzig.de

Abstract

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26

Knowledge of plant species distributions is essential for various applications, such as nature conservation, agriculture, and forestry. Remote sensing data, especially high-resolution orthoimages from Unoccupied Aerial Vehicles (UAVs), were demonstrated to be an effective data source for plant species mapping. Particularly, in concert with novel pattern recognition methods, such as Convolutional Neural Networks (CNNs), plant species can be accurately segmented in such high-resolution UAV images. Training such pattern recognition models for species segmentation that are transferable across various landscapes and remote sensing data characteristics often requires excessive training data. Training data are usually derived in the form of segmentation masks from field surveys or visual interpretation of the target species in remote sensing images. Still, both methods are laborious and constrain the training of transferable pattern recognition models. Alternatively, pattern recognition models could be trained on the open knowledge of how plants look as available from smartphone-based species identification apps, that is, millions of citizen science-based smartphone photographs and the corresponding species label. However, these pairs of citizen science-based photographs and simple species labels (one label for the entire image) cannot be used directly for training state-of-the-art segmentation models used for UAV image analysis, which require per-pixel labels for training (also called masks). Here, we overcome the limitation of simple labels of citizen science plant observations with a two-step approach: In the first step, we train CNN-based image classification models using the simple labels and apply them in a moving-window approach over UAV orthoimagery to create segmentation masks. In the second phase, these segmentation masks are used to train state-of-the-art CNN-based image segmentation models with an encoder-decoder structure. We tested the approach on UAV orthoimages acquired in summer and autumn on a test site comprising ten temperate deciduous tree species in varying mixtures. Several tree species could be mapped with surprising accuracy (mean



27 F1-score = 0.47). In homogenous species assemblages, the accuracy increased considerably
28 (mean F1-score 0.55). The results indicate that many tree species can be mapped without
29 generating training data and by integrating pre-existing knowledge from citizen science.
30 Moreover, our analysis revealed that citizen science photographs' variability in acquisition
31 data and context facilitates the generation of models that are transferable through the
32 vegetation season. Thus, citizen science data may greatly advance our capacity to monitor
33 hundreds of plant species and, thus, Earth's biodiversity across space and time.

34 Keywords: Remote Sensing, Convolutional Neural Network, Citizen Science Data,
35 Plant species, Transfer learning.

36 1 Introduction

37 Spatially explicit information on plant species is crucial for various applications, including na-
38 ture conservation, agriculture, and forestry. Remote sensing emerged as a promising tool to
39 create spatially continuous maps of plant species (Müllerová et al., 2023; Bouguettaya et al.,
40 2022; Fassnacht et al., 2016). Thereby, supervised machine learning algorithms are commonly
41 used to identify species-specific features in spatial, temporal, or spectral patterns of remotely
42 sensed signals (Sun et al., 2021; Maes and Steppe, 2019; Lopatin et al., 2019; Curnick et al.,
43 2021; Wagner, 2021). In recent years, remote sensing imagery from drones, also known as
44 Unoccupied Air Vehicles (UAVs), has emerged as an effective source of information for map-
45 ping plant species (Kattenborn et al., 2021; Fassnacht et al., 2016; Schiefer et al., 2020). By
46 means of mosaicing a series of individual image frames, UAVs enable the creation of georef-
47 erenced orthoimagery of relatively large areas with extremely high spatial resolution, e.g., in
48 the mili- or centimeter range. The fine spatial grain of such imagery can reveal distinctive
49 morphological plant features to identify specific plant species. Such plant features include
50 the leaf shape, flowers, branching patterns, or crown structures (Sun et al., 2021; Kattenborn
51 et al., 2019a). An effective way to unleash the potential of these fine spatial features is given
52 by deep learning-based pattern-recognition techniques, in particular by Convolutional Neural
53 Networks (CNN). A series of studies have demonstrated that CNN can precisely segment plant
54 species' crowns in high-resolution UAV imagery (Kattenborn et al., 2021; Hoeser and Kuen-
55 zer, 2020; Brodrick et al., 2019). Such CNN models learn the characteristic spatial features
56 of the target (here, plant species) through a cascade of filter operations (convolutions). Given
57 these high-dimensional computations, efficiently adopting these models to UAV orthoimagery,
58 with their large spatial extents but also high resolution, requires training and applying them
59 sequentially using smaller sub-regions of an orthoimage (e.g., image tiles of 512 by 512 pixels,
60 Fig. 1a).

61 However, generating models that are transferable across various landscapes and remote
62 sensing data characteristics requires large amounts of training data (Kattenborn et al., 2021;
63 Galuszynski et al., 2022). In particular, when neighboring plant species bear a similar resem-
64 blance, a wealth of training data becomes essential, allowing the model to discern the subtle
65 distinctions between these species (Kattenborn et al., 2021; Schiefer et al., 2020). Commonly,
66 the generation of training data is costly. Training data are usually derived from field surveys



67 or visual interpretation of remote sensing images, also known as annotation or labelling. Both
68 methods have limitations: Field surveys are often logistically challenged by site accessibility
69 or travel costs. Moreover, field surveys commonly only enable the acquisition of point ob-
70 servations or relative cover fractions of the target species (Leitão et al., 2018). Visual image
71 interpretation is often much more effective (Kattenborn et al., 2019b; Schiefer et al., 2023)
72 but for some species, precise visual identification of species can be challenging due to sub-
73 tle indicative morphological features, the variability of these features in the landscape, or
74 the complexity of vegetation communities (e.g., smooth transitions of canopies of different
75 species). Moreover, the representativeness of data derived from field surveys and visual in-
76 terpretation is often limited to the location where and when the data were acquired, which
77 may reduce a model’s generalization to new regions or time periods (Kattenborn et al., 2022).
78 Therefore, the obtained amount and quality of training data can be a critical factor for the
79 performance and transferability of CNN models (Bayraktar et al., 2020; Rzanny et al., 2019;
80 Brandt et al., 2020).

81 The challenge of limited training data for UAV-based plant species identification may
82 be alleviated by the collective power of scientists and citizens openly sharing their plant
83 observations on the web (Ivanova and Shashkov, 2021; Fraisl et al., 2022; Di Cecco et al.,
84 2021). A particular data treasure in this regard is generated by citizen science projects
85 for plant species identification. Examples are the iNaturalist and Pl@ntNet projects, which
86 encourage ten-thousands of individuals to capture, share, and annotate photographs of the
87 World’s plant life (Boone and Basille, 2019; Di Cecco et al., 2021). The quantity of such
88 citizen science observations is rapidly growing due to the increasing number of volunteers
89 participating in the platform (Boone and Basille, 2019; Di Cecco et al., 2021).

90 Currently, the iNaturalist project contains over 26 Mio of globally distributed and anno-
91 tated photographs of vascular plant species. The iNaturalist platform allows users to identify
92 plant species manually or using a computer vision model integrated into the platform. The
93 submitted observations are then evaluated by the community, and a research-grade classifica-
94 tion is assigned if over two-thirds of the community agrees on the species identification. The
95 Pl@ntNet project includes over 20 Mio observations of globally distributed vascular plants.
96 Pl@ntNet requires users to photograph their observations and select an organ tag (e.g., leaf,
97 flower, fruit, or stem). The Pl@ntNet features an image recognition algorithm to analyze
98 the tagged photograph and suggest a plant species. Pl@ntNet’s validation process uses a
99 dynamic approach, combining automated algorithm confidence with community consensus
100 (Joly et al., 2016). The validated observations of iNaturalist and Pl@ntNet are shared via the
101 Global Biodiversity Information Facility (GBIF), a global network that provides open access
102 to biodiversity data (GBIF, 2019).

103 Citizen science-based plant photographs with species annotations provide a valuable, large,
104 and continuously growing data source for training pattern recognition models, such as CNNs
105 (Van Horn et al., 2018; Joly et al., 2016). However, such citizen science data has a cardinal
106 limitation: It only provides simple species annotation for a plant photograph (*the image_i*
107 *shows species_j*). Hence, these labels only enable to train image classification models that



108 predict the likelihood of a species being present in an image but not where in the image.
109 In an ideal setting for species mapping applications, the species labels would delineate the
110 regions or pixels belonging to a species (*The pixels in the right corner of image_i represents*
111 *species_j*). Such labels (known as masks) could be used to train CNN-based segmentation
112 models, which can predict a species probability for each individual pixel of an image (or tile
113 of an orthoimage) (Galuszynski et al., 2022; Schiefer et al., 2020).

114 In a pioneering study by Soltani et al. (2022), the limitation of the simple labels that
115 come with citizen science photographs was overcome by a workaround. At first, image classi-
116 fication models were trained with citizen science data and simple labels to predict a species
117 per image. The trained image classification models were then applied sequentially on tiles
118 of 512×512 pixels of UAV-based orthomosaics in a moving-window-like fashion with very
119 high overlap (Fig. 1a). Lastly, the individual predictions derived from the moving-window
120 steps were rasterized to a seamless segmentation map (Fig. 1b). However, this workaround
121 is computationally intense and inefficient for large or multiple UAV orthomosaics, as seg-
122 mentation maps can only be derived from many overlapping prediction steps. In contrast,
123 state-of-the-art CNN-based segmentation methods (typically an encoder-decoder structure)
124 used in remote sensing applications are trained with reference data in the form of masks with
125 dimensions (pixels) corresponding to the extent of the imagery, where each pixel of the mask
126 defines the absence or presence of a class (here plant species) in the imagery (Kattenborn
127 et al., 2021). Respective segmentation models are more efficient as they segment multiple
128 classes in a single prediction step. Moreover, they enable more detailed class representations
129 in situations where multiple classes are arranged in complex patterns.

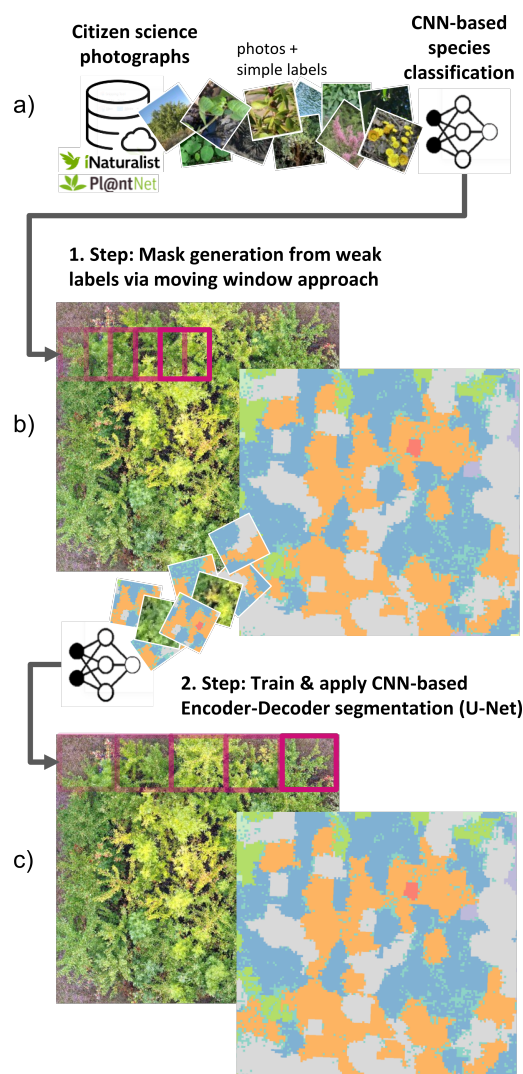


Figure 1: 1-column figure: Schematic representation of the proposed workflow, including the moving window approach by [Soltani et al. \(2022\)](#) (a,b) and the use of state-of-the-art encoder-decoder segmentation algorithms (c).

130 Here, we propose a solution to overcome the limitation of simple annotations of citizen
131 science plant observations with a two-step approach: In the first step, we apply the procedure
132 of [Soltani et al. \(2022\)](#), involving CNN-based image classification models trained on citizen
133 science photographs and simple species labels to predict plant species in UAV orthoimages
134 using the moving-window approach described above (Fig. 1a, b). Although computationally
135 demanding, this serves to create segmentation masks for UAV orthoimages. In the second step,
136 these segmentation masks are used to train more efficient CNN-based image segmentation
137 models with an encoder-decoder structure (Fig. 1c). These more efficient models could then
138 be applied to larger spatial extents or due new UAV orthomosaics (e.g. of different sites or



139 time steps).

140 The present study, hence, addresses the following research questions:

- 141 • Can we harness weak labels from citizen science plant observations to train efficient
142 state-of-the-art semantic segmentation models?
- 143 • Do those segmentation models also increase the accuracy compared to the simple moving
144 window approach?

145 These questions are evaluated on a tree species dataset acquired on an experimental site
146 (MyDiv experiment, Bad Lauchstädt, Germany), where ten temperate deciduous tree species
147 were planted in stratified and complex mixtures. The selection of this location is attributed
148 to its harmonious coexistence of various plant species within a compact area.

149 2 Methods

150 2.1 Data acquisition and pre-processing

151 2.1.1 Study site and drone data acquisition

152 The MyDiv experimental site is located in Bad Lauchstädt, Saxony-Anhalt, Germany (lati-
153 tude, 51°23' N, longitude, 11°53' E). The site comprises 80 plots composed in different configu-
154 rations of ten deciduous tree species, including *Acer pseudoplatanus*, *Aesculus hippocastanum*,
155 *Betula pendula*, *Carpinus betulus*, *Fagus sylvatica*, *Fraxinus excelsior*, *Prunus avium*, *Quercus*
156 *petraea*, *Sorbus aucuparia*, and *Tilia platyphyllos* (Ferlian et al., 2018). Each plot measures
157 12 m by 12 m and contains 140 trees planted at distances of 1 m (Fig 2). In total, all plots
158 together accommodate 11,200 individual trees. Each plot contains varying tree species com-
159 positions, including one, two, and four tree species. This variety in species, their balanced
160 composition, and plots of different canopy complexity (species mixtures) provide an ideal
161 setting to test the proposed species segmentation approach.

162 We collected UAV-based RGB aerial imagery over the MyDiv experimental site using a
163 DJI Mavic 2 Pro and the flight planning software DroneDeploy (DroneDeploy vers. 5.0, USA).
164 Two flights were conducted in 2022 in July and September, where July corresponds to the peak
165 of the growing season and September to senescence stage (Fig 2). The flight plan was setup
166 with a forward overlap of 90%, side overlap of 70% at an altitude of 16 m (ground sampling
167 distance of approximately 0.22 cm per pixel). We used the generated images and Metashape
168 (vers. 1.7.6, Agisoft LLC) to generate orthoimages for both flight campaigns. The orthoimage
169 for July and September are onward called Ortho_{July} and Ortho_{September}, respectively.

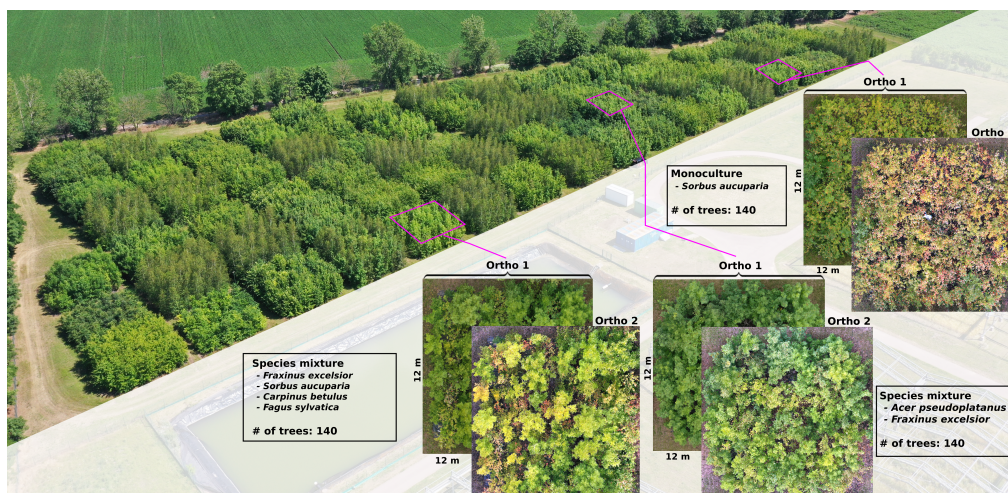


Figure 2: Overview of the MyDiv experimental site with close-ups for three plots of different species composition. The MyDiv site is located at Lat. 51.3916 N, Long. 11.8857 E.

170 To evaluate the performance of the CNN models for tree species mapping, we created
171 reference data by manually delineating the tree species in the UAV orthoimages in QGIS
172 (vers. 3.32.3). To reduce the workload, we did not delineate the species for the entire plot
173 but for diagonal transects with 20 m length and 2 m width.

174 2.1.2 Citizen science training data

175 We queried plant observations of the iNaturalist and Pl@ntNet projects via the GBIF database
176 for our target tree species using scientific names. For the iNaturalist data, we used the
177 R package rinat (vers. 0.1.8), an API to iNaturalist. The Pl@ntNet data were acquired
178 by submitting a download request for the selected tree species via GBIF. The number of
179 photographs available from iNaturalist and Pl@ntNet varied for the different tree species.
180 Per species, we were able to acquire between 582 to 10000 photographs (mean 7696) from
181 the iNaturalist platform and 221 to 3304 images (mean 2238) from the Pl@ntNet platform
182 (details see Appendix Table A1).

183 In addition to the tree species, we added a background class to consider canopy gaps
184 between trees. For training data, we used the Google Image API to query different keywords,
185 e.g. *grass*, *forest floor*, *forest ground*. After cleaning the obtained images for non-meaningful
186 results, the background class included 1100 photographs.

187 We converted all photographs to a rectangular shape by cropping them to the shorter side
188 and resampled them to a common size of 512×512 pixels (the tile size used later for the CNN
189 model generation). Figure 3 shows examples of the downloaded photographs for the different
190 tree species and a comparison to their appearance in Ortho_{July}.

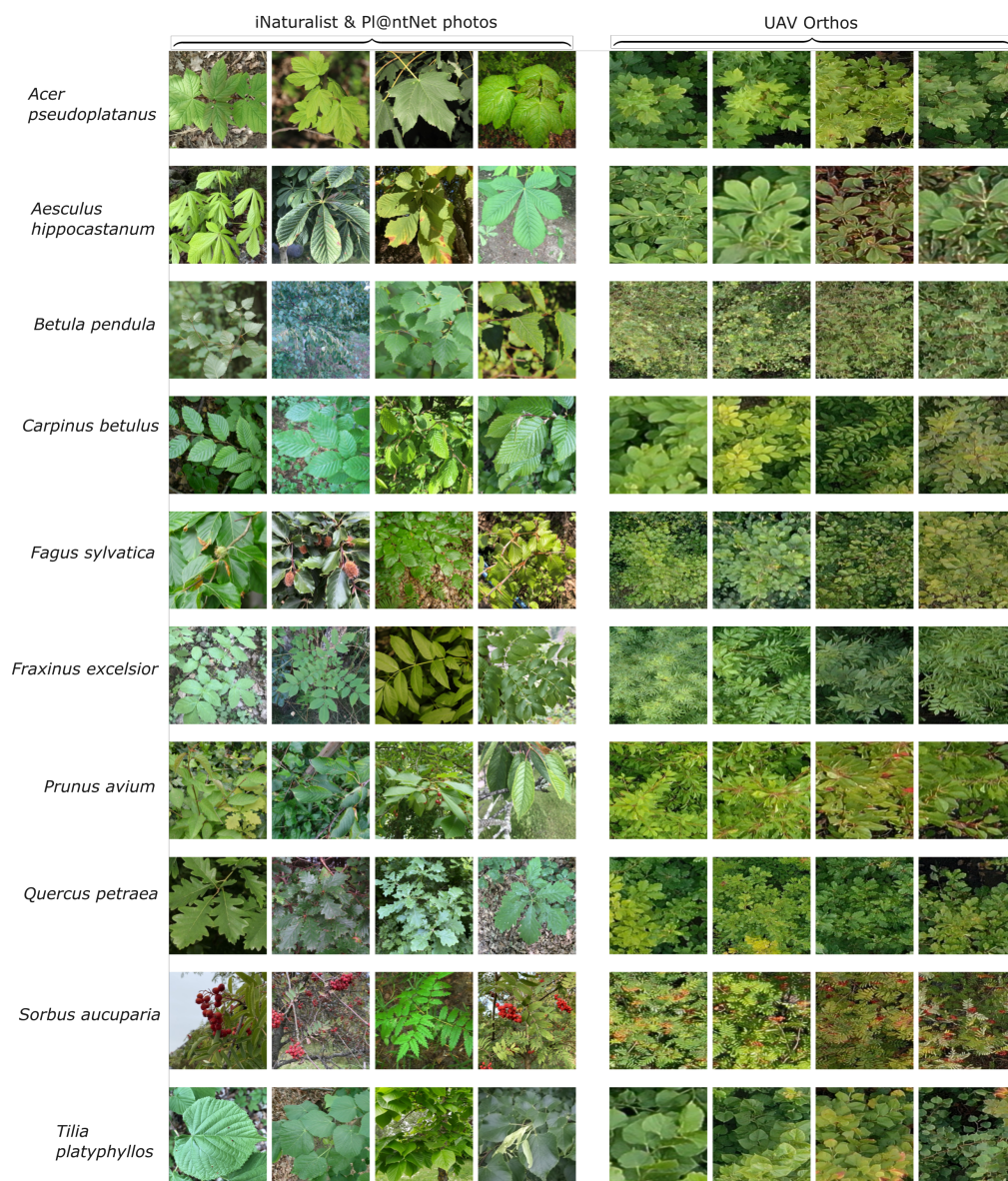


Figure 3: Example citizen science-based photographs derived from iNaturalist and tiles of UAV orthoimages (512 * 512 pixels) for the ten tree species in the MyDiv experiment.

191 The acquisition settings of citizen science plant photographs are heterogeneous and differ
192 considerably from the typical bird perspective of UAV orthoimages. For instance, from the
193 UAV perspective, canopies are mostly viewed from a relatively homogeneous distance, and
194 the photographs represent mostly leaves and other crown components. In contrast, the citi-
195 zen science data includes a lot of close-ups, landscape imagery, or horizontal photographs of



196 trunks. [Soltani et al. \(2022\)](#) has demonstrated that species recognition in UAV images can be
197 improved by excluding crowd-sourced photographs that are exceptionally close (e.g., showing
198 individual leaf veins) or too far away from the plant (e.g., landscape images). Accordingly,
199 we filtered the citizen science-based training photos according to the camera-plant-distance.
200 Moreover, we filtered photos that exclusively contained tree stems. Because such information
201 is unavailable in the citizen science datasets, we trained CNN-based regression and classifi-
202 cation models to predict acquisition distance and tree trunk presence for each downloaded
203 photograph. To train these CNN-based models, we visually estimated the acquisition distance
204 (4,500 photographs) and labeled tree trunk presence (1,000 photographs). To ease the label-
205 ing process, we used previously labeled training data from ([Soltani et al., 2022](#)) and added
206 150 additional tree photographs from the tree species present in the MyDiv experimental site.

207 To predict acquisition distance and trunk presence, We randomly split the citizen science-
208 based plant photographs into training and validation sets, with 80% for training and 20% for
209 validation.

210 For the distance regression and the trunk classification, we used the EfficientNetB7 back-
211 bone ([Tan and Le, 2019](#)). For the distance regression, we used the following top-layer settings:
212 global average pooling, batch normalization, drop out (rate 0.1), and a final dense layer with
213 1 unit and linear activation function. We used the Adam optimizer (learning rate of 0.0001)
214 and a mean squared error (MSE) loss function. For the trunk classification, we used the
215 following top-layer settings: global max-pooling, a final dense layer with two units, and a
216 softmax activation function. We used the Adam optimizer (learning rate of 0.0001) and the
217 categorical cross-entropy loss function. Both models were trained using a batch size of 20 and
218 50 epochs.

219 We used the model with the lowest loss from these epochs (details on the model perfor-
220 mance are given in Appendix [A1.3](#)) to predict the acquisition distance and tree trunk presence
221 in all downloaded photographs for our target species. We filtered training photographs prior
222 to training CNN-based species classification (see section [2.2](#)) with acquisition distances less
223 than 0.2 m and greater than 15 m and photographs classified as trunk (probability threshold
224 of 0.5). Thereby, 82,628 of the 101,574 downloaded citizen science photographs remained.

225 **2.2 CNN-based creation of plant species segmentation masks using a mov-** 226 **ing window approach**

227 The segmentation masks were obtained using a CNN image classification model trained on
228 crowd-sourced plant photographs and simple species labels using a moving window method
229 (hereafter $\text{CNN}_{\text{window}}$, Fig. 1). Based on the results of previous studies, we choose a generic
230 image size of 512×512 pixels for the CNN classification model ([Schiefer et al., 2020](#); [Soltani](#)
231 [et al., 2022](#)). During the moving window approach, the orthoimage is sequentially cropped
232 into tiles of 512×512 pixels on which the image classification is applied to predict the species
233 for each location. This procedure is applied with a dense overlap between tiles defined by
234 a step size, resulting in a dense regular grid of species predictions. We chose a vertical and
235 horizontal distance of 51 pixels as step size. The resulting predictions are afterward rasterized



236 to a continuous species distribution grid with a spatial resolution of 8.31 cm/pixel (see [Soltani](#)
237 [et al., 2022](#), for details). The $\text{CNN}_{\text{window}}$ model was implemented as a classification task with
238 eleven classes, including the ten tree species and the background class.

239 The number of available photographs varied widely across tree species (see [2.1.2](#)), poten-
240 tially biasing the model towards classes with more photographs. To address this imbalance,
241 we equally sampled 4,000 photographs for each class with replacement. We applied a data
242 augmentation to increase the variance of the duplicated images. The augmentation consisted
243 of random vertical and horizontal flips, random brightness maximum delta of 10% (± 0.1), and
244 contrast alteration within a range of 90% to 110% (0.9 to 1.1) of training photographs. We
245 randomly partitioned the training data into validation and training sets to ensure unbiased
246 evaluation. We allocated a holdout of 20% of the training data for model selection, while the
247 remaining 80% was used for model training. Subsequently, we assessed the accuracy of the
248 selected model using independent reference data.

249 After testing different architectures as model backbones, including ResNet-50V2, Effi-
250 cientNetB07, and EfficientNetV2L, we selected EfficientNetV2L. The following layers were
251 added on top of the EfficientNetV2L backbone: Dropout with a ratio of 0.5, average pooling,
252 dropout with a ratio of 0.5, dense layer with 128 units, L2 kernel regularizer (0.001), a ReLU
253 activation function, and a final dense layer with a softmax activation function and 11 units.
254 We used Root Mean Squared Propagation (RMSprop) as the optimizer with a learning rate
255 of 0.0001 and categorical cross-entropy as a loss function. We trained the configured model
256 with a batch size of 15 over 150 epochs. The model with the lowest validation loss (based
257 on the 20% holdout) was selected as the final model. The latter was used to predict the tree
258 species (probabilities) in the UAV orthoimages using the abovementioned $\text{CNN}_{\text{window}}$ method.
259 To filter uncertain predictions (predominantly in canopy gaps or at crown shadows), we only
260 considered a tree species as predicted above a threshold higher than 0.6. Otherwise, it was
261 assigned to NA (not available). To smooth the predictions and remove noise, we applied
262 a sieve operation on the output of the $\text{CNN}_{\text{window}}$ (threshold = 50, considering horizontal,
263 vertical, and diagonal neighbors, R-package *terra*, vers. 1.7).

264 **2.3 CNN-based plant species segmentation using an encoder-decoder ar-** 265 **chitecture**

266 As encoder-decoder segmentation architecture (onwards $\text{CNN}_{\text{segment}}$), we chose U-Net ([Ron-](#)
267 [neberger et al., 2015](#)), which is the most widely applied segmentation method in remote sensing
268 image segmentation ([Kattenborn et al., 2021](#)). The U-Net architecture is a CNN-based algo-
269 rithm that performs semantic segmentation by predicting a class for each pixel of the input
270 image. The architecture consists of an encoder-decoder structure with skip connections. The
271 configured architecture has four levels of convolutional blocks. Each convolutional block con-
272 sists of two convolutional layers and is followed by batch normalization and ReLU activation.
273 The encoder gradually compresses feature maps and reduces their spatial dimensions via max
274 pooling operations, while the decoder increases the feature map resolution by transposed con-
275 volution. The encoder and decoder blocks are connected through skip connections, which



276 transfer the spatial context of the encoder feature maps to the decoder, enabling a segmen-
277 tation at high-resolution in the last layer. The final layer has eleven units (corresponding to
278 the ten tree species and a background class). A corresponding softmax activation function
279 maps the features to class probabilities. Using a max function, the pixels of the segmentation
280 output are assigned to the class with the highest probability (Fig. A12).

281 The segmentation masks for training $\text{CNN}_{\text{segment}}$ were obtained from the predictions of the
282 $\text{CNN}_{\text{window}}$ method applied on both UAV orthoimages (section 2.2, $\text{Ortho}_{\text{July}}$, $\text{Ortho}_{\text{September}}$).
283 At first, we resampled the $\text{CNN}_{\text{window}}$ prediction maps to the original spatial resolution of the
284 orthoimages (0.22 cm pixel size). Afterward, we cropped the orthoimages and the prediction
285 maps into non-overlapping tiles, each with a size of 512×512 pixels, resulting in a total of
286 44,980 and 37,113 tiles from $\text{Ortho}_{\text{July}}$ and $\text{Ortho}_{\text{September}}$, respectively.

287 The training data obtained from the $\text{CNN}_{\text{window}}$ approach were filtered to avoid training
288 the $\text{CNN}_{\text{segment}}$ with uncertain predictions. Thereby, we assumed that higher model uncer-
289 tainty are present in areas where the model predicts multiple classes with low relative cover.
290 Thus, after initial tests, we included only those tiles where the cover of at least one class
291 exceeded 30%. The number of training tiles per class after filtering varied between 1257 and
292 16894 samples; *Acer pseudoplatanus* (6581), *Aesculus hippocastanum* (2054), *Betula pendula*
293 (4955), *Carpinus betulus* (1535), *Fagus sylvatica* (16894), *Fraxinus excelsior* (7901), *Prunus*
294 *avium* (1257), *Quercus petraea* (1302), *Sorbus aucuparia* (5473), *Tilia platyphyllos* (1982),
295 Background (5408).

296 Similar to the previous $\text{CNN}_{\text{window}}$ classification task, the availability of training tiles
297 varied greatly across the tree species. This class imbalance may have partially stemmed from
298 the more systematic misclassification of certain classes during the $\text{CNN}_{\text{window}}$ prediction. To
299 reduce the unfavorable effects of a class imbalance on model training, we sampled 4,000 tiles
300 per class with replacements (similar to the $\text{CNN}_{\text{window}}$ procedure). We applied the same
301 data augmentation strategy as $\text{CNN}_{\text{window}}$ to increase variance among duplicates. 20% of the
302 training data were withheld for model selection.

303 We trained the U-Net architecture using Root Mean Squared Propagation (RMSprop) as
304 the optimizer with a learning rate of 0.0001 and an adapted Dice loss function. We adapted
305 the Dice loss to ignore the weights coming from pixels with NA mask values. The models
306 were trained with a batch size of 20 over 150 epochs.

307 The $\text{CNN}_{\text{segment}}$ was then applied to $\text{Ortho}_{\text{July}}$ and $\text{Ortho}_{\text{September}}$. To reduce uncertain
308 predictions of $\text{CNN}_{\text{segment}}$, we assigned the pixels where predicted probabilities did not exceed
309 0.3 to the background class. Thereby, we assumed that uncertain predictions predominantly
310 occur in canopy gaps. As image segmentations typically suffer from increased uncertainty at
311 tile edges, we repeated the predictions with horizontal and vertical shifts of 256 pixels, which
312 were subsequently aggregated using a majority vote.

313 The final model performance of $\text{CNN}_{\text{segment}}$ was assessed and compared to $\text{CNN}_{\text{window}}$
314 using the independent reference data (transects) obtained from the visual interpretation of
315 the UAV orthoimages.



316 3 Results

317 For the $\text{CNN}_{\text{window}}$ method, F1-scores differed considerably across the tree species, while
318 these differences were relatively consistent across the two orthoimages, i.e. $\text{Ortho}_{\text{July}}$ and
319 $\text{Ortho}_{\text{September}}$ (Fig. 4a, b). On a plot level, comparably high model performance (mean F1 >
320 0.6) was found for *Acer pseudoplatanus* and *Fraxinus excelsior*, followed by the intermediate
321 performance (mean F1-score 0.35-0.55) for *Aesculus hippocastanum*, *Sorbus aucuparia*, *Tilia*
322 *platyphyllos*, *Betula pendula*, and *Carpinus betulus*. Low performance (mean F1-score < 0.35)
323 was found for *Quercus petraea*, *Fagus sylvatica*, and *Prunus avium*. Averaged across species,
324 there was a slight decrease in model performance from $\text{Ortho}_{\text{July}}$ with a mean F1-score of 0.44
325 to $\text{Ortho}_{\text{September}}$ with a mean F1-score of 0.4 (Fig. 4a, b). Note that $\text{Ortho}_{\text{July}}$ corresponded
326 to the peak of the season, where leaves and canopies were still fully developed.

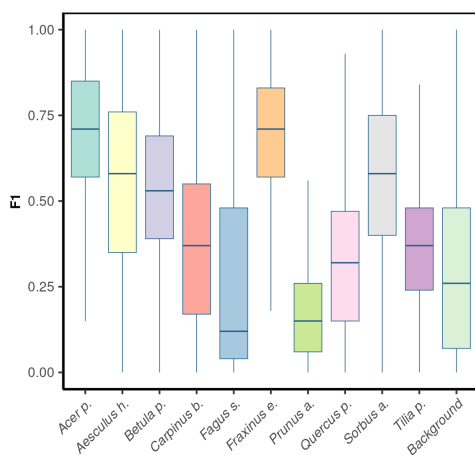
327 The $\text{CNN}_{\text{segment}}$ model performance across species was similar but generally higher com-
328 pared to the $\text{CNN}_{\text{window}}$ method. For $\text{Ortho}_{\text{July}}$ F1-scores increased from 0.44 to 0.48 (Fig. 4a
329 vs. c) and for $\text{Ortho}_{\text{September}}$, F1-scores increased from 0.40 to 0.46 (Fig. 4b vs. d).

330 We observed notable differences in model performance (mean F1) across different species
331 mixtures, which are plots having one, two, or four species per plot (Fig. 5). For both
332 $\text{CNN}_{\text{window}}$ and $\text{CNN}_{\text{segment}}$, the model performance strongly increased with lower number
333 of species per plot (results for $\text{CNN}_{\text{window}}$ are given in the Appendix; Fig. A13).

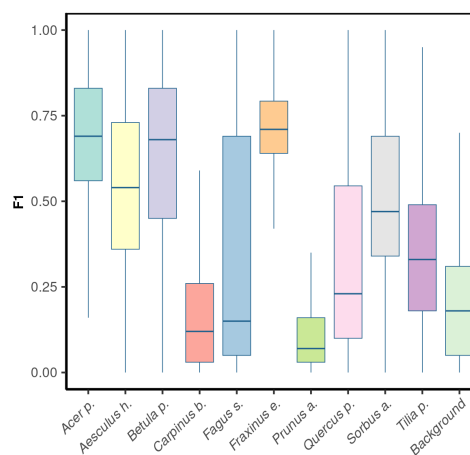
334 The model performance of $\text{CNN}_{\text{segment}}$ exceeded the model performance of $\text{CNN}_{\text{window}}$
335 particularly in plots with increased number of species: For monocultures the relative increase
336 in model performance (F1-score) amounted to 2.5%, in two species plots to 6.9%, and in
337 plots with four species to 20.9% (averaged for $\text{Ortho}_{\text{July}}$ and $\text{Ortho}_{\text{September}}$). This increased
338 performance can be attributed to the advantages of the encoder-decoder principle of the
339 $\text{CNN}_{\text{segment}}$ method, enabling a pixel-wise and contextual prediction at the original resolution
340 of the orthomosaics. These advantages are also visible in Fig. 6, where $\text{CNN}_{\text{segment}}$ resulted
341 in more detailed and accurate tree species segmentations (particularly for plot 26 and 29).

342 The highest model performance for $\text{CNN}_{\text{segment}}$ was found in monoculture plots, where
343 F1-scores > 0.5 was found for eight out ten species for both $\text{Ortho}_{\text{July}}$ and $\text{Ortho}_{\text{September}}$. A
344 considerably lower performance for the July and September acquisition was found for *Prunus*
345 *avium*, which may correspond to similarities in leaf and canopy structure with *Fagus sylvatica*
346 and *Fraxinus excelsior* (a confusion matrix is given in the Appendix, Fig. A11). The decreased
347 performance for *Carpinus betulus* and *Prunus avium* in $\text{Ortho}_{\text{September}}$ can be attributed to
348 the very advanced senescence and leaf loss.

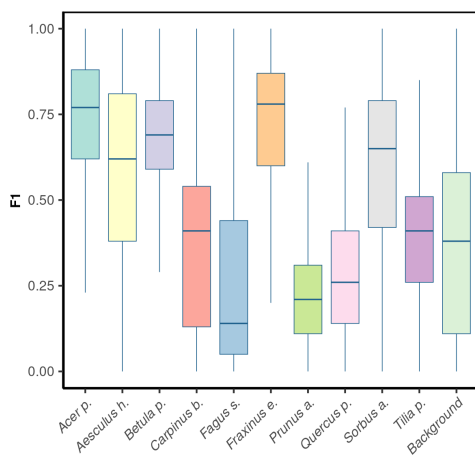
349 In addition to the increase in model performance, our analysis revealed that the prediction
350 on orthoimagery using $\text{CNN}_{\text{segment}}$ only required 10% of the computation time compared to
351 $\text{CNN}_{\text{window}}$. The duration of applying the models to the whole MyDiv orthomosaics covering
352 an area of (3.02 hectare; 0.22 cm ground sampling distance) took approximately 27.05 hours
353 with $\text{CNN}_{\text{segment}}$ and 264.88 hours with $\text{CNN}_{\text{window}}$ (NVIDIA A6000 with 48 GB RAM).



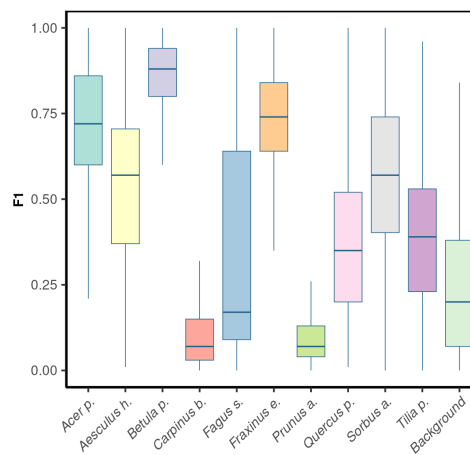
(a) F1-scores for CNN_{window} on OrthoJuly (mean 0.44).



(b) F1-scores of CNN_{window} on OrthoSeptember (mean 0.42).

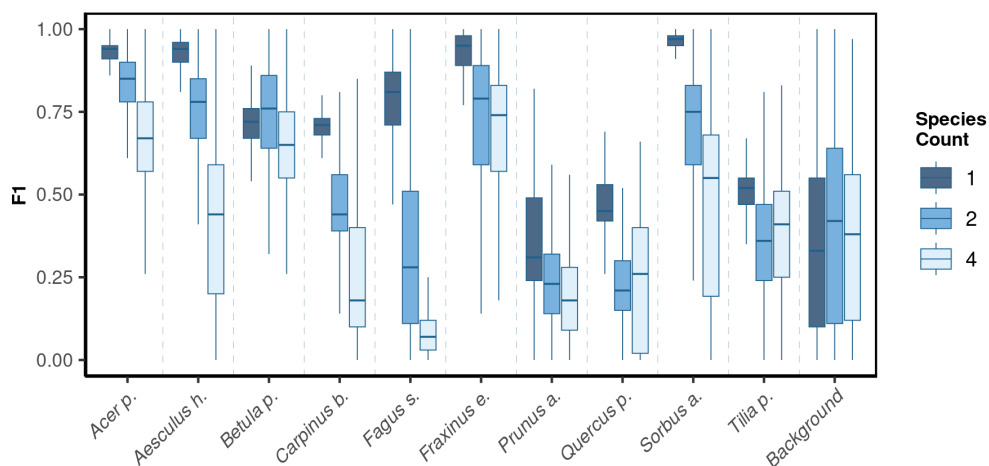


(c) F1-scores of $CNN_{segment}$ on OrthoJuly (mean 0.48).

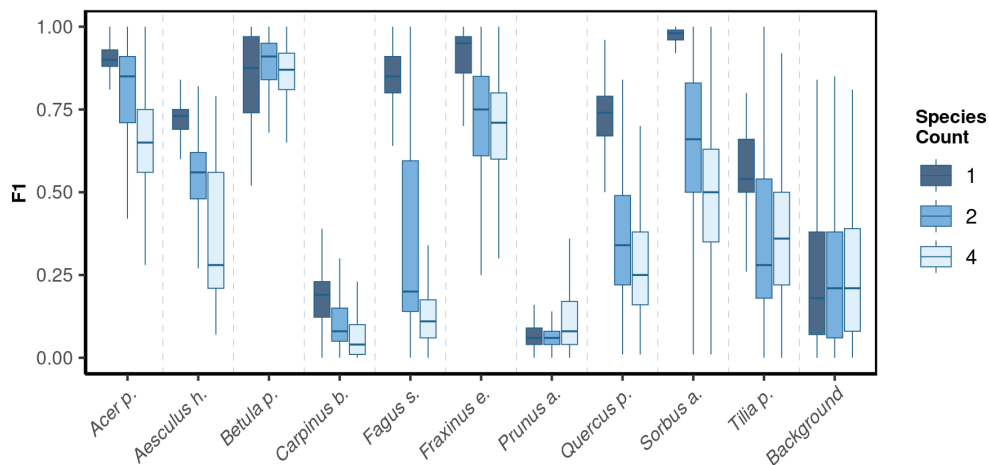


(d) F1-scores of $CNN_{segment}$ on OrthoSeptember (mean 0.46).

Figure 4: F1-scores by tree species and background class for OrthoJuly and OrthoSeptember derived from CNN_{window} and $CNN_{segment}$.



(a) Performance across species mixtures (F1-scores) on Ortho_{July}. Mean F1-scores: 1 species (0.51), 2 species (0.44), 4 species (0.41)



(b) Performance across species mixtures (F1-scores) on Ortho_{September}. Mean F1-scores: 1 species (0.58), 2 species (0.51), 4 species (0.42)

Figure 5: The model performance (F1-score) of the CNN_{segment} model across a gradient of canopy complexity in Ortho_{July} and Ortho_{September}. F1-scores decrease with increasing canopy complexity in plots



Figure 6: Transects of 2 m by 20 m of selected plots, including the orthoimage, the reference, CNN_{window} predictions, and CNN_{segment} predictions. Visualizations for the remaining plots are given in the Appendix (Section A1.1).



354 4 Discussion

355 4.1 Filtering of citizen science data for drone-related applications

356 To achieve better correspondence between plant features visible in the citizen science pho-
357 tographs and the UAV images, we filtered the crowd-sourced photographs based on their
358 acquisition distance (less than 0.3 m or greater than 15 m) to exclude macro and landscape
359 photographs. Moreover, we excluded photographs that predominantly display tree trunks,
360 facilitating a foliage-centric perspective as intrinsic to high-resolution UAV images (Fig. 3).
361 In the future, more criteria may be used for filtering citizen science imagery, including meta-
362 data (labels) on the presence of specific plant organs within an image (e.g., fruits, flowers) as
363 provided as a by-product by some citizen science plant identification apps (e.g., Pl@ntNet).

364 4.2 The creation of segmentation masks from simple image labels

365 One of the challenges of generating segmentation masks for the encoder-decoder method
366 ($\text{CNN}_{\text{segment}}$) with the proposed workflow may be error propagation between the different
367 steps. Firstly, the CNN image classification trained on the citizen science data has varyin
368 uncertainty for the different species, resulting from noisy citizen science observations or lim-
369 itations to identify some species solely by photographs (Van Horn et al., 2018). Secondly,
370 the moving window approach ($\text{CNN}_{\text{window}}$), which predicts one species for an entire tile, may
371 be too coarse to resemble very complex canopies (e.g., in highly diverse plant communities).
372 However, although the fact that the segmentation labels created with the $\text{CNN}_{\text{window}}$ approach
373 are partially relatively inaccurate (Fig 4a, 6), we found that the $\text{CNN}_{\text{segment}}$ procedure in-
374 deed resulted in higher performance than the $\text{CNN}_{\text{window}}$ procedure. This is in line with
375 other studies (Kattenborn et al., 2021; Cloutier et al., 2023; Schiller et al., 2021) reporting
376 that deep learning-based pattern recognition can partially overcome noise labels, whereas the
377 intentional use of noisy reference data, also known as weakly-supervised learning, is gener-
378 ally very promising in the absence of high-quality labels (Zhou, 2018). Here, we filtered the
379 training data (masks) for regions where we expect extreme noise levels, that is, for tiles where
380 none of the classes exceeded a relative cover of 30%. These regions were, according to our
381 observation, often canopy gaps and shadowed areas, where one naturally expects lower model
382 performance due to less distinct species-specific textures (Lopatin et al., 2019; Milas et al.,
383 2017; De Sa et al., 2018).

384 The enhanced segmentation performance of the $\text{CNN}_{\text{segment}}$ approach compared to $\text{CNN}_{\text{window}}$
385 can be attributed to the spatially explicit and finer-resolved predictions of the U-Net segmen-
386 tation algorithm (encoder-decoder principle), enabling to segment the tree species at the
387 native resolution of the orthoimagery. Particularly, for plots with more species (two or four)
388 the encoder-decoder segmentation approach resulted in improved prediction results compared
389 to the $\text{CNN}_{\text{window}}$ method in plots with more species (two or four) and hence, more complex
390 canopies. Thus, the presented two-step approach of creating segmentation masks from sim-
391 ple class labels $\text{CNN}_{\text{window}}$, as provided by iNaturalist and Pl@ntNet platforms, can indeed
392 be used to create segmentation masks required for state-of-the-art image analysis methods



393 (CNN_{segment}) and thereby result in higher value for remote sensing applications. The in-
394 creased value of these segmentation masks enables the training of algorithms with higher
395 performance in species recognition. It greatly enhances the efficiency of applying the models
396 on orthoimagery (factor of approximately ten). Especially for recurrent applications, such
397 as monitoring or large-scale undertakings, the two-step approach involving the creation of
398 segmentation masks and encoder-decoder architectures is recommended.

399 4.3 The role of canopy complexity

400 Overall, the segmentation performance declined with increasing species richness per plot.
401 We expect that this can mainly be attributed to the small size of individual trees at the
402 MyDiv site, where in high species mixtures, there is a lower chance that a 512×512 pixel
403 tile includes clearly visible species-specific leaf and branching patterns. This also explains
404 why, in particular, trees with lower relative canopy height (e.g., *Quercus petrea* and *Fagus*
405 *sylvatica*) were less likely to be accurately segmented in species mixtures. The observed effect
406 of canopy complexity is in line with previous findings from [Soltani et al. \(2022\)](#); [Lopatin](#)
407 [et al. \(2017\)](#); [Fassnacht et al. \(2016\)](#); [Fricker et al. \(2019\)](#), where smaller patches of individual
408 species were less likely to be accurately detected. Visual inspection also confirmed that
409 false predictions were more likely at canopy edges between different tree species (Fig. 6).
410 However, it should be noted that the small-scaled canopy complexity of the plots used here
411 is exceptionally high (Fig. 3). Most tree crowns in the MyDiv experiment do not exceed a
412 diameter of 1.5 m, and the transition among tree crowns of multiple species is often very
413 fuzzy. Thus, we expect reduced performance in canopy transitions to be less relevant in
414 real-world settings, where tree species appear in more extensive, homogeneous patches and
415 where individual crowns are commonly larger. Thus, the model performance in these species
416 mixtures can be interpreted as a rather conservative estimate. The results obtained for the
417 monocultures might be more representative in terms of real-world applications, as mature
418 trees in temperate forests typically have crown diameters 5 to 20 times larger. Application
419 tests of the presented approach in real forests are desirable. However, acquiring such a dataset
420 is a logistical challenge since temperate forest stands commonly do not feature a comparably
421 high and balanced occurrence of that many tree species.

422 4.4 Spatial resolution of the UAV imagery is key

423 According to the results obtained in the monocultures, The CNN_{segment} model successfully
424 classified seven out of ten tree species (F1 > 0.7). The lower F1-scores for *Quercus petrea*
425 (mean F1 0.57), *Prunus avium*(mean F1 0.2), *Tilia platyphyllos*(mean F1 0.53) may result
426 from the spectral and morphological similarity at the current spatial resolution of the UAV
427 imagery (0.22 cm)(Fig. 3). Hence, there was a tendency that these species were often confused
428 with each other (see confusion matrices in Appendix A1.2). Such confusion among plants
429 with a similar appearance was confirmed by other studies ([Cloutier et al., 2023](#); [Schiefer](#)
430 [et al., 2020](#), e.g.) and matches our experience from the generation of reference data via visual



431 interpretation, where a separation between these species was sometimes challenging. Initial
432 CNN-based segmentation attempts (results not shown) in the preparation of this study were
433 based on an orthoimage of 0.3 cm instead of 0.22 cm resolution, resulting in clearly lower
434 model performances. This aligns with the reported importance of spatial resolution of UAV
435 imagery for CNN segmentation of earlier studies (Schiefer et al., 2020; Schmitt et al., 2020;
436 Ma et al., 2019; G. Braga et al., 2020). Thus, while the current orthoimages with 0.22 cm
437 resolution delivered promising results, further increasing the spatial resolution might be very
438 promising for species where characteristic leaf forms can only be visualized at fine spatial
439 resolutions.

440 4.5 Model transferability across seasons and orthoimage acquisition prop- 441 erties

442 The diversity of human behavior and electronic devices makes citizen science-based plant
443 photographs very heterogeneous. This can be a challenge for deep learning applications, such
444 as species recognition or plant trait characterization (Schiller et al., 2021; Van Horn et al.,
445 2021; van Der Velde et al., 2023; Affouard et al., 2017), where models have to identify features
446 that hold across various viewing angles, distances, or illumination conditions. However, this
447 heterogeneity might also be of great value, given that citizens depict the appearance of plants
448 under various site, environmental, and phenological conditions. This, in turn, offers a unique
449 setting for training models that are generic and transferable across these conditions. Here, we
450 evaluated the transferability of our models across different data sets by applying them to two
451 orthoimages acquired in different seasons (peak of growing season and autumn). Both the
452 CNN_{window} and CNN_{segment} models could identify deciduous tree species in the orthoimages
453 with surprising accuracies, suggesting that the models are transferable to different conditions.

454 4.6 Outlook

455 Overall, our results indeed highlight the value of citizen science photographs with simple
456 class labels to create training data for state-of-the-art segmentation approaches. A great
457 advantage of this citizen science-based approach is that it does not require commonly costly
458 training data obtained from visual interpretation or field surveys (here, we only acquired
459 reference data for validating the procedure). This particularly highlights the potential of
460 citizen science data for applications where many species are of interest, such as biodiversity-
461 related monitoring applications (Chandler et al., 2017; Johnston et al., 2023). In this regard,
462 data or models of species-recognition platforms that incorporate excessive amounts of plant
463 species and respective imagery are very promising, including iNaturalist (Boone and Basille,
464 2019), Pl@ntNet (Affouard et al., 2017), ObsIdentify (Molls, 2021) or FloraIncognita (Mäder
465 et al., 2021). Yet, based on the current and the precursor study (Soltani et al., 2022), we
466 expect that a pre-selection of citizen science photograph databases considering images more
467 representative of the common UAV-based perspective is required to unleash the potential of
468 this heterogeneous data.



469 5 Conclusion

470 The transfer learning approach presented here demonstrates the value of freely available
471 crowd-sourced plant photographs for remote sensing studies. This heterogeneous dataset
472 can provide valuable training data for transferable CNN-based segmentation models. Here,
473 this potential was highlighted in a very complex task, i.e., the differentiation of multiple tem-
474 perate deciduous tree species in mixed vegetation stands with a complex structure. The
475 presented two-step approach demonstrated how we can transfer and harness generic knowl-
476 edge gathered by citizens on how plants 'look' to the bird perspective of high-resolution drone
477 imagery. The presented moving window approach overcomes the limitation of citizen science-
478 based photographs having only simple species labels. The segmentation maps derived from
479 an image classification model applied in a moving window setting can be harnessed to create
480 segmentation masks for encoder-decoder-type segmentation models. The latter does not only
481 enable higher accuracies in species segmentation but is also considerably more efficient. By
482 building on the effort of thousands of citizens, this framework enables the mapping of plant
483 species without any training data obtained from visual interpretation or ground-based field
484 surveys. Due to the excessive amounts of plant photographs acquired in different conditions,
485 such models can be assumed to have a large transferability.

486 6 Data and code availability

487 The code used in this study is publicly accessible via our GitHub repository at [https://](https://github.com/salimsoltani28/CrowdVision2TreeSegment)
488 github.com/salimsoltani28/CrowdVision2TreeSegment. The data supporting the findings
489 of this research is available on Zonodo at <https://zenodo.org/uploads/10019552>.

490 7 Declaration of competing interest

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

493 8 Acknowledgements

494 SS and TK acknowledge funding by the German Research Foundation (DFG) under the
495 project BigPlantSens (Assessing the Synergies of Big Data and Deep Learning for the Re-
496 mote Sensing of Plant Species; Project number 444524904) and PANOPS (Revealing Earth's
497 plant functional diversity with citizen science; project number 504978936). SS and HF ac-
498 knowledge financial support by the Federal Ministry of Education and Research of Germany
499 (BMBF) and by the Saechsische Staatsministerium für Wissenschaft, Kultur und Tourismus
500 in the program Center of Excellence for AI-research "Center for Scalable Data Analytics and
501 Artificial Intelligence Dresden/Leipzig", project identification number: ScaDS.AI.NE and OF
502 acknowledge funding by the Deutsche Forschungsgemeinschaft DFG (German Centre for In-



503 tegrative Biodiversity Research, FZT118; and Gottfried Wilhelm Leibniz Prize, Ei 862/29-1).
504 Moreover, we acknowledge support from Leipzig University for Open Access Publishing.

505 References

- 506 A. Affouard, H. Goëau, P. Bonnet, J.-C. Lombardo, and A. Joly. Pl@ ntnet app in the era of
507 deep learning. In *ICLR: International Conference on Learning Representations*, 2017.
- 508 E. Bayraktar, M. E. Basarkan, and N. Celebi. A low-cost uav framework towards ornamental
509 plant detection and counting in the wild. *ISPRS Journal of Photogrammetry and Remote*
510 *Sensing*, 167:1–11, 2020.
- 511 M. E. Boone and M. Basille. Using inaturalist to contribute your nature observations to
512 science. *EDIS*, 2019(4):5–5, 2019.
- 513 A. Bouguettaya, H. Zarzour, A. Kechida, and A. M. Taberkit. Deep learning techniques to
514 classify agricultural crops through uav imagery: A review. *Neural Computing and Applica-*
515 *tions*, 34(12):9511–9536, 2022.
- 516 M. Brandt, C. J. Tucker, A. Kariryaa, K. Rasmussen, C. Abel, J. Small, J. Chave, L. V.
517 Rasmussen, P. Hiernaux, A. A. Diouf, et al. An unexpectedly large count of trees in the
518 west african sahara and sahel. *Nature*, 587(7832):78–82, 2020.
- 519 P. G. Brodrick, A. B. Davies, and G. P. Asner. Uncovering ecological patterns with convolu-
520 tional neural networks. *Trends in ecology & evolution*, 34(8):734–745, 2019.
- 521 M. Chandler, L. See, K. Copas, A. M. Bonde, B. C. López, F. Danielsen, J. K. Legind,
522 S. Masinde, A. J. Miller-Rushing, G. Newman, et al. Contribution of citizen science towards
523 international biodiversity monitoring. *Biological conservation*, 213:280–294, 2017.
- 524 M. Cloutier, M. Germain, and E. Laliberté. Influence of temperate forest autumn leaf phenol-
525 ogy on segmentation of tree species from uav imagery using deep learning. *bioRxiv*, pages
526 2023–08, 2023.
- 527 D. J. Curnick, A. J. Davies, C. Duncan, R. Freeman, D. M. Jacoby, H. T. Shelley, C. Rossi,
528 O. R. Wearn, M. J. Williamson, and N. Pettorelli. Smallsats: a new technological frontier
529 in ecology and conservation? *Remote Sensing in Ecology and Conservation*, 2021.
- 530 N. C. De Sa, P. Castro, S. Carvalho, E. Marchante, F. A. López-Núñez, and H. Marchante.
531 Mapping the flowering of an invasive plant using unmanned aerial vehicles: is there potential
532 for biocontrol monitoring? *Frontiers in plant science*, 9:293, 2018.
- 533 G. J. Di Cecco, V. Barve, M. W. Belitz, B. J. Stucky, R. P. Guralnick, and A. H. Hurlbert.
534 Observing the observers: How participants contribute data to inaturalist and implications
535 for biodiversity science. *BioScience*, 71(11):1179–1188, 2021.



- 536 F. E. Fassnacht, H. Latifi, K. Stereńczak, A. Modzelewska, M. Lefsky, L. T. Waser, C. Straub,
537 and A. Ghosh. Review of studies on tree species classification from remotely sensed data.
538 *Remote Sensing of Environment*, 186:64–87, 2016.
- 539 O. Ferlian, S. Cesarz, D. Craven, J. Hines, K. E. Barry, H. Bruelheide, F. Buscot, S. Haider,
540 H. Heklau, S. Herrmann, et al. Mycorrhiza in tree diversity–ecosystem function relation-
541 ships: conceptual framework and experimental implementation. *Ecosphere*, 9(5):e02226,
542 2018.
- 543 D. Fraisl, G. Hager, B. Bedessem, M. Gold, P.-Y. Hsing, F. Danielsen, C. B. Hitchcock, J. M.
544 Hulbert, J. Piera, H. Spiers, et al. Citizen science in environmental and ecological sciences.
545 *Nature Reviews Methods Primers*, 2(1):64, 2022.
- 546 G. A. Fricker, J. D. Ventura, J. A. Wolf, M. P. North, F. W. Davis, and J. Franklin. A
547 convolutional neural network classifier identifies tree species in mixed-conifer forest from
548 hyperspectral imagery. *Remote Sensing*, 11(19):2326, 2019.
- 549 J. R. G. Braga, V. Peripato, R. Dalagnol, M. P. Ferreira, Y. Tarabalka, L. E. OC Aragão,
550 H. F. de Campos Velho, E. H. Shiguemori, and F. H. Wagner. Tree crown delineation
551 algorithm based on a convolutional neural network. *Remote Sensing*, 12(8):1288, 2020.
- 552 N. C. Galuszynski, R. Duker, A. J. Potts, and T. Kattenborn. Automated mapping of por-
553 tulacaria afro canopies for restoration monitoring with convolutional neural networks and
554 heterogeneous unmanned aerial vehicle imagery. *PeerJ*, 10:e14219, 2022.
- 555 GBIF. Gbif: the global biodiversity information facility, 2019.
- 556 T. Hoeser and C. Kuenzer. Object detection and image segmentation with deep learning on
557 earth observation data: A review-part i: Evolution and recent trends. *Remote Sensing*, 12
558 (10):1667, 2020.
- 559 N. Ivanova and M. Shashkov. The possibilities of gbif data use in ecological research. *Russian*
560 *Journal of Ecology*, 52:1–8, 2021.
- 561 A. Johnston, E. Matechou, and E. B. Dennis. Outstanding challenges and future directions
562 for biodiversity monitoring using citizen science data. *Methods in Ecology and Evolution*,
563 14(1):103–116, 2023.
- 564 A. Joly, P. Bonnet, H. Goëau, J. Barbe, S. Selmi, J. Champ, S. Dufour-Kowalski, A. Affouard,
565 J. Carré, J.-F. Molino, et al. A look inside the pl@ ntnet experience: The good, the bias
566 and the hope. *Multimedia Systems*, 22:751–766, 2016.
- 567 T. Kattenborn, J. Eichel, and F. E. Fassnacht. Convolutional neural networks enable effi-
568 cient, accurate and fine-grained segmentation of plant species and communities from high-
569 resolution uav imagery. *Scientific reports*, 9(1):1–9, 2019a.



- 570 T. Kattenborn, J. Lopatin, M. Förster, A. C. Braun, and F. E. Fassnacht. Uav data as
571 alternative to field sampling to map woody invasive species based on combined sentinel-1
572 and sentinel-2 data. *Remote sensing of environment*, 227:61–73, 2019b.
- 573 T. Kattenborn, J. Leitloff, F. Schiefer, and S. Hinz. Review on convolutional neural networks
574 (cnn) in vegetation remote sensing. *ISPRS Journal of Photogrammetry and Remote Sensing*,
575 173:24–49, 2021.
- 576 T. Kattenborn, F. Schiefer, J. Frey, H. Feilhauer, M. D. Mahecha, and C. F. Dormann.
577 Spatially autocorrelated training and validation samples inflate performance assessment
578 of convolutional neural networks. *ISPRS Open Journal of Photogrammetry and Remote
579 Sensing*, 5:100018, 2022.
- 580 P. J. Leitão, M. Schwieder, F. Pötzschner, J. R. Pinto, A. M. Teixeira, F. Pedroni, M. Sanchez,
581 C. Rogass, S. van der Linden, M. M. Bustamante, et al. From sample to pixel: multi-scale
582 remote sensing data for upscaling aboveground carbon data in heterogeneous landscapes.
583 *Ecosphere*, 9(8):e02298, 2018.
- 584 J. Lopatin, F. E. Fassnacht, T. Kattenborn, and S. Schmidtlein. Mapping plant species in
585 mixed grassland communities using close range imaging spectroscopy. *Remote Sensing of
586 Environment*, 201:12–23, 2017.
- 587 J. Lopatin, K. Dolos, T. Kattenborn, and F. E. Fassnacht. How canopy shadow affects invasive
588 plant species classification in high spatial resolution remote sensing. *Remote Sensing in
589 Ecology and Conservation*, 5(4):302–317, 2019.
- 590 L. Ma, Y. Liu, X. Zhang, Y. Ye, G. Yin, and B. A. Johnson. Deep learning in remote sensing
591 applications: A meta-analysis and review. *ISPRS journal of photogrammetry and remote
592 sensing*, 152:166–177, 2019.
- 593 P. Mäder, D. Boho, M. Rzanny, M. Seeland, H. C. Wittich, A. Deggelmann, and J. Wäldchen.
594 The flora incognita app—interactive plant species identification. *Methods in Ecology and
595 Evolution*, 2021.
- 596 W. H. Maes and K. Steppe. Perspectives for remote sensing with unmanned aerial vehicles
597 in precision agriculture. *Trends in plant science*, 24(2):152–164, 2019.
- 598 A. S. Milas, K. Arend, C. Mayer, M. A. Simonson, and S. Mackey. Different colours of
599 shadows: Classification of uav images. *International Journal of Remote Sensing*, 38(8-10):
600 3084–3100, 2017.
- 601 C. Molls. The obs-services and their potentials for biodiversity data assessments with a test
602 of the current reliability of photo-identification of coleoptera in the field. *Tijdschrift voor
603 Entomologie*, 164(1-3):143–153, 2021.



- 604 J. Müllerová, G. Brundu, A. Große-Stoltenberg, T. Kattenborn, and D. M. Richardson. Pat-
605 tern to process, research to practice: remote sensing of plant invasions. *Biological Invasions*,
606 pages 1–26, 2023.
- 607 O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image
608 segmentation. In *International Conference on Medical image computing and computer-*
609 *assisted intervention*, pages 234–241. Springer, 2015.
- 610 M. Rzanny, P. Mäder, A. Deggelmann, M. Chen, and J. Wäldchen. Flowers, leaves or both?
611 how to obtain suitable images for automated plant identification. *Plant Methods*, 15(1):
612 1–11, 2019.
- 613 F. Schiefer, T. Kattenborn, A. Frick, J. Frey, P. Schall, B. Koch, and S. Schmidlein. Mapping
614 forest tree species in high resolution uav-based rgb-imagery by means of convolutional neural
615 networks. *ISPRS Journal of Photogrammetry and Remote Sensing*, 170:205–215, 2020.
- 616 F. Schiefer, S. Schmidlein, A. Frick, J. Frey, R. Klinke, K. Zielewska-Büttner, S. Junttila,
617 A. Uhl, and T. Kattenborn. Uav-based reference data for the prediction of fractional cover
618 of standing deadwood from sentinel time series. *ISPRS Open Journal of Photogrammetry*
619 *and Remote Sensing*, 8:100034, 2023.
- 620 C. Schiller, S. Schmidlein, C. Boonman, A. Moreno-Martínez, and T. Kattenborn. Deep
621 learning and citizen science enable automated plant trait predictions from photographs.
622 *Scientific Reports*, 11(1):1–12, 2021.
- 623 M. Schmitt, J. Prexl, P. Ebel, L. Liebel, and X. X. Zhu. Weakly supervised semantic seg-
624 mentation of satellite images for land cover mapping—challenges and opportunities. *arXiv*
625 *preprint arXiv:2002.08254*, 2020.
- 626 S. Soltani, H. Feilhauer, R. Duker, and T. Kattenborn. Transfer learning from citizen science
627 photographs enables plant species identification in uavs imagery. *ISPRS Open Journal of*
628 *Photogrammetry and Remote Sensing*, page 100016, 2022.
- 629 Z. Sun, X. Wang, Z. Wang, L. Yang, Y. Xie, and Y. Huang. Uavs as remote sensing platforms
630 in plant ecology: review of applications and challenges. *Journal of Plant Ecology*, 14(6):
631 1003–1023, 2021.
- 632 M. Tan and Q. Le. Efficientnet: Rethinking model scaling for convolutional neural networks.
633 In *International conference on machine learning*, pages 6105–6114. PMLR, 2019.
- 634 M. van Der Velde, H. Goëau, P. Bonnet, R. d’Andrimont, M. Yordanov, A. Affouard,
635 M. Claverie, B. Czúcz, N. Elvekjær, L. Martinez-Sanchez, et al. Pl@ ntnet crops: merg-
636 ing citizen science observations and structured survey data to improve crop recognition for
637 agri-food-environment applications. *Environmental Research Letters*, 18(2):025005, 2023.



- 638 G. Van Horn, O. Mac Aodha, Y. Song, Y. Cui, C. Sun, A. Shepard, H. Adam, P. Perona, and
639 S. Belongie. The inaturalist species classification and detection dataset. In *Proceedings of*
640 *the IEEE conference on computer vision and pattern recognition*, pages 8769–8778, 2018.
- 641 G. Van Horn, E. Cole, S. Beery, K. Wilber, S. Belongie, and O. Mac Aodha. Benchmarking
642 representation learning for natural world image collections. In *Proceedings of the IEEE/CVF*
643 *Conference on Computer Vision and Pattern Recognition*, pages 12884–12893, 2021.
- 644 F. H. Wagner. The flowering of atlantic forest pleroma trees. *Scientific reports*, 11(1):1–20,
645 2021.
- 646 Z.-H. Zhou. A brief introduction to weakly supervised learning. *National science review*, 5
647 (1):44–53, 2018.



648 **A Appendix**

649 **A1.1 Prediction maps**

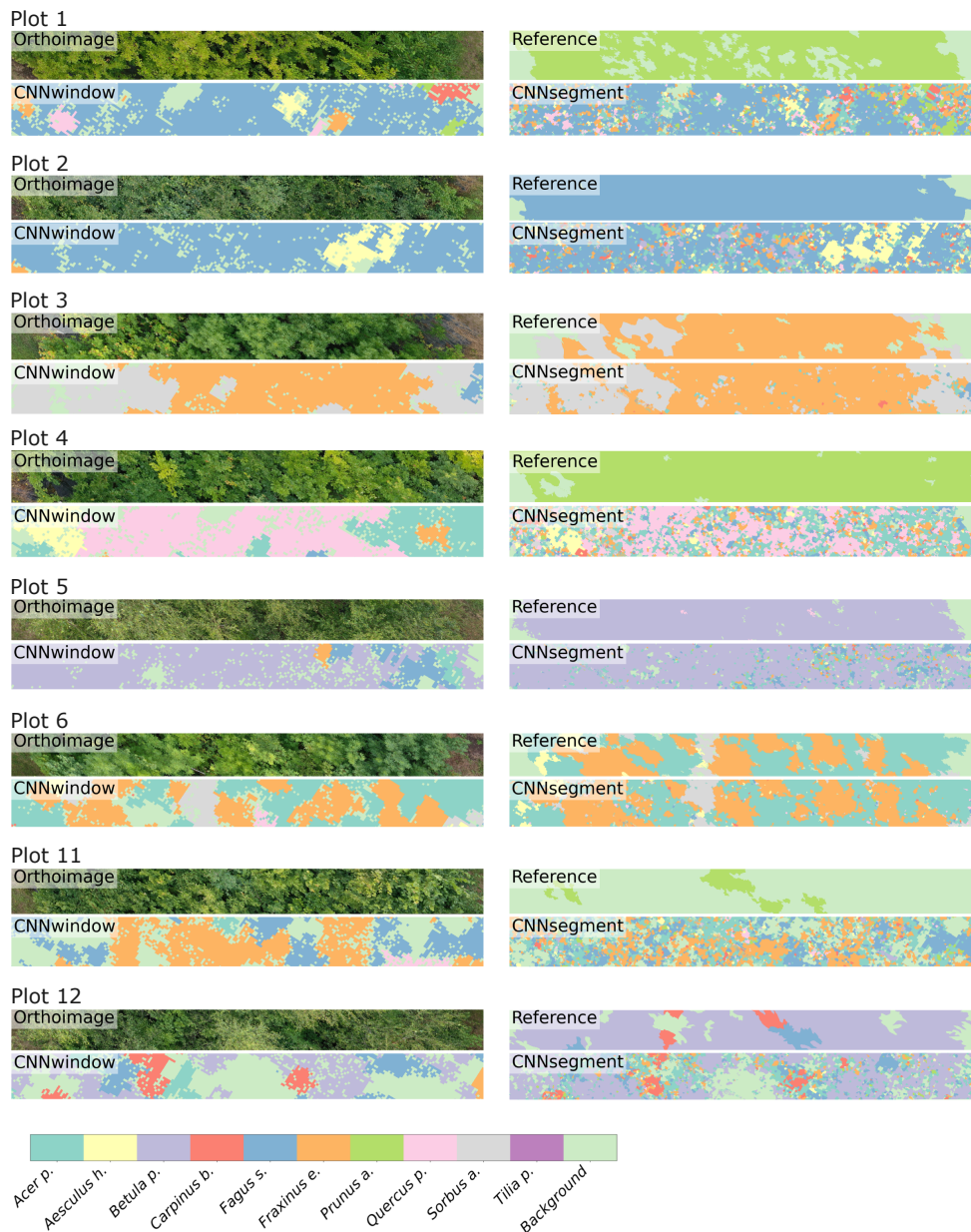


Figure A1: Transects of 2 m by 20 m of selected plots, including the orthoimage, the reference, CNN_{window} predictions, and CNN_{segment} predictions.



Figure A2: Transects of 2 m by 20 m of selected plots, including the orthoimage, the reference, CNN_{window} predictions, and CNN_{segment} predictions.

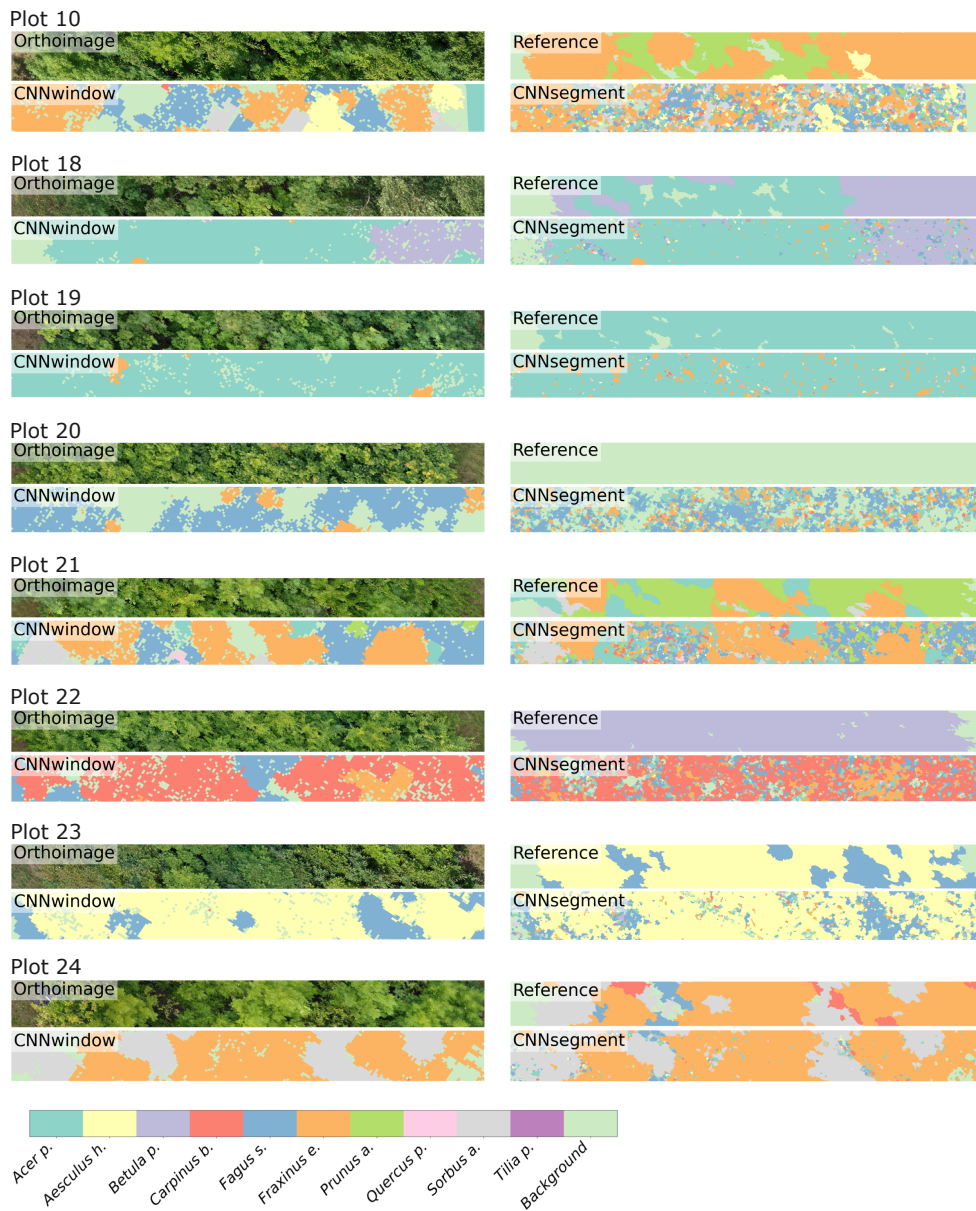


Figure A3: Transects of 2 m by 20 m of selected plots, including the orthoimage, the reference, CNN_{window} predictions, and CNN_{segment} predictions.

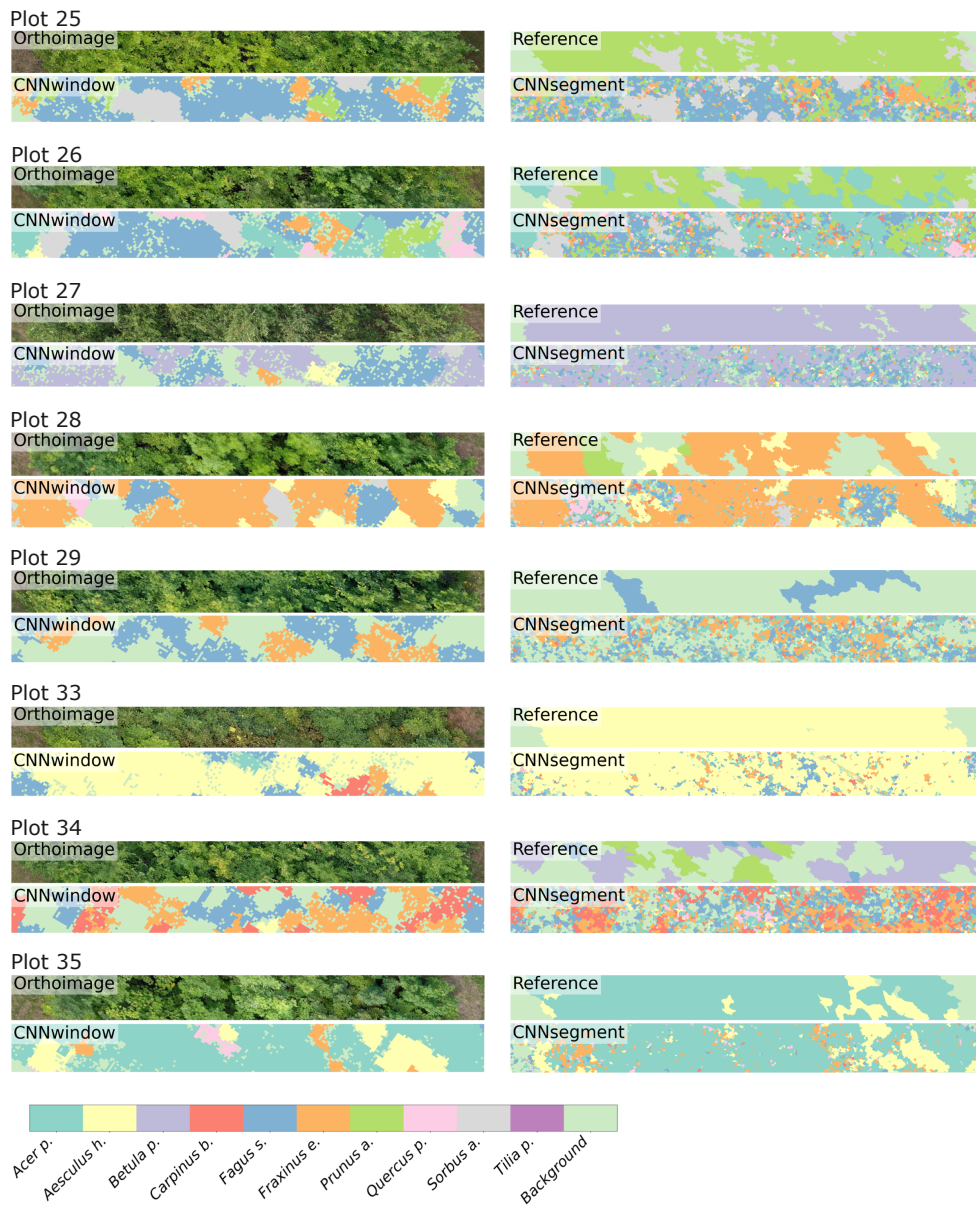


Figure A4: Transects of 2 m by 20 m of selected plots, including the orthoimage, the reference, CNN_{window} predictions, and CNN_{segment} predictions.

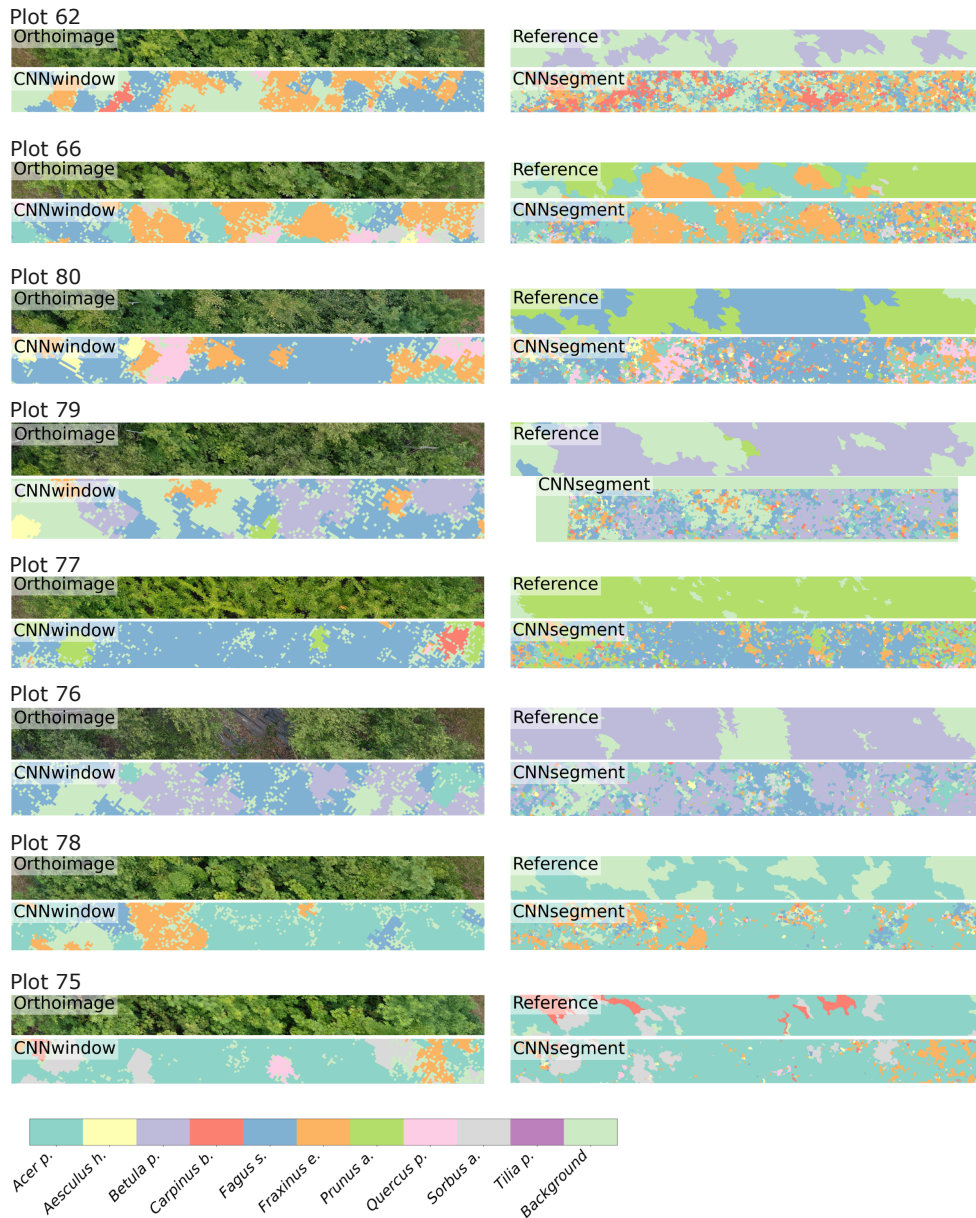


Figure A5: Transects of 2 m by 20 m of selected plots, including the orthoimage, the reference, CNN_{window} predictions, and CNN_{segment} predictions.

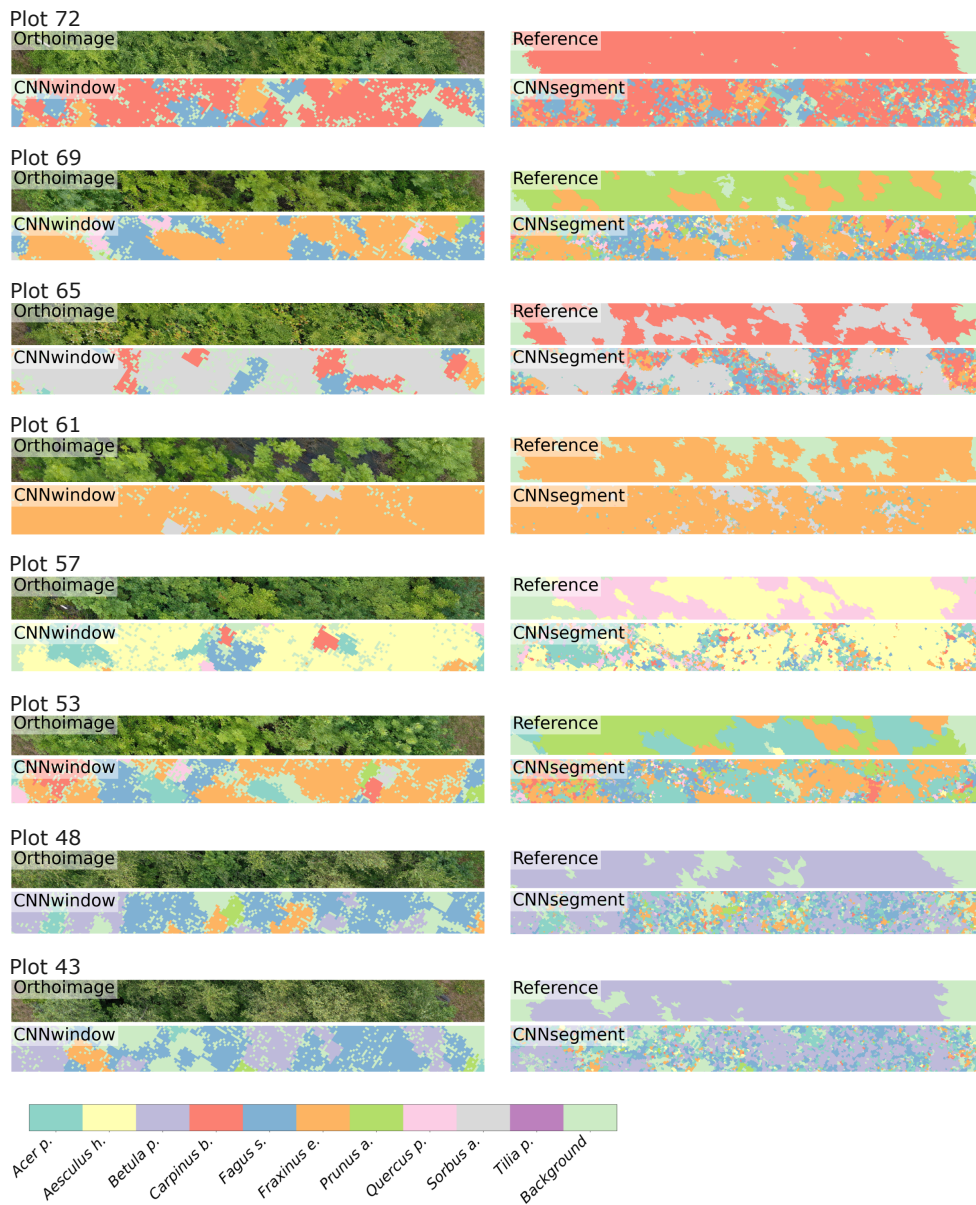


Figure A6: Transects of 2 m by 20 m of selected plots, including the orthoimage, the reference, CNN_{window} predictions, and CNN_{segment} predictions.

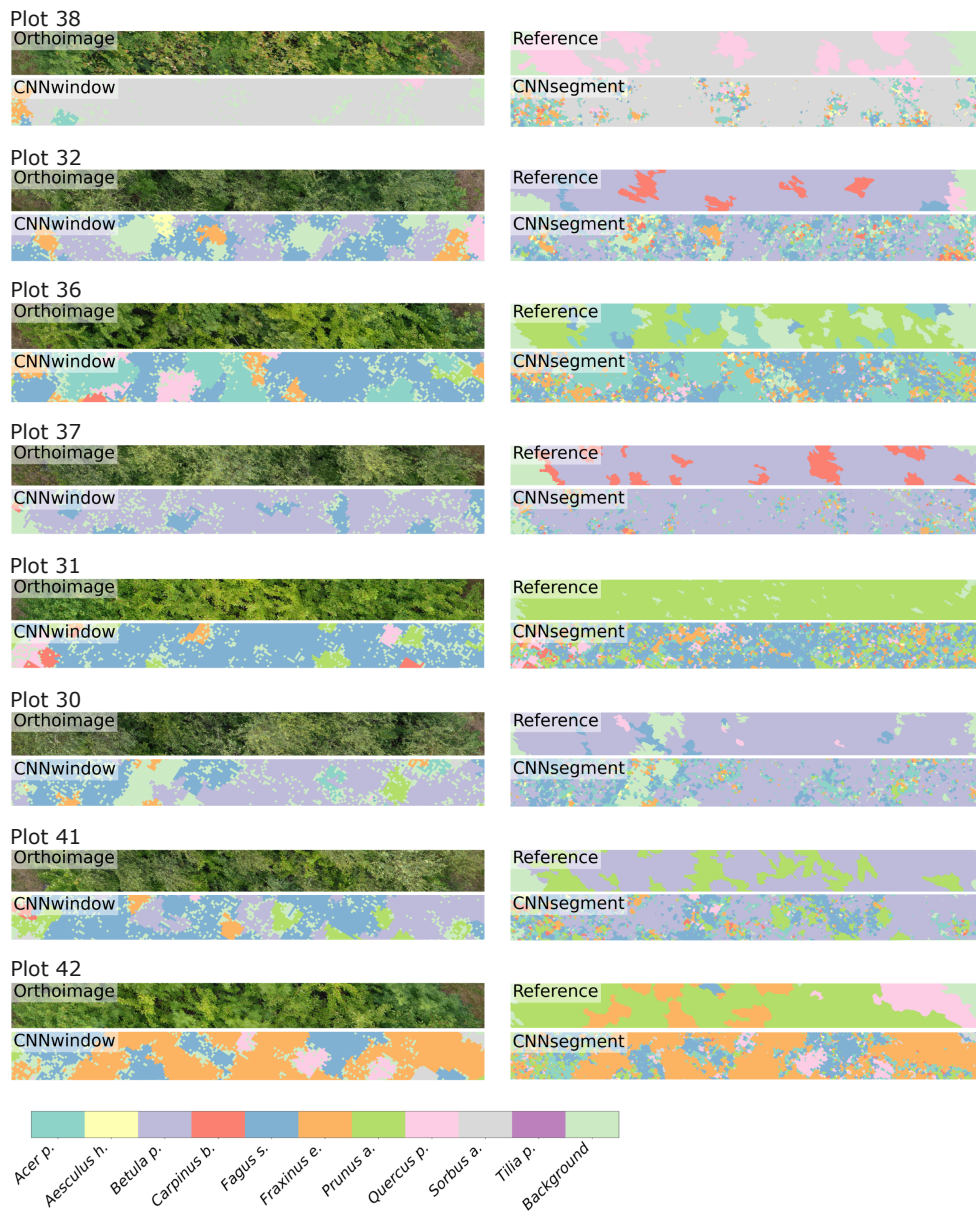


Figure A7: Transects of 2 m by 20 m of selected plots, including the orthoimage, the reference, CNN_{window} predictions, and CNN_{segment} predictions.

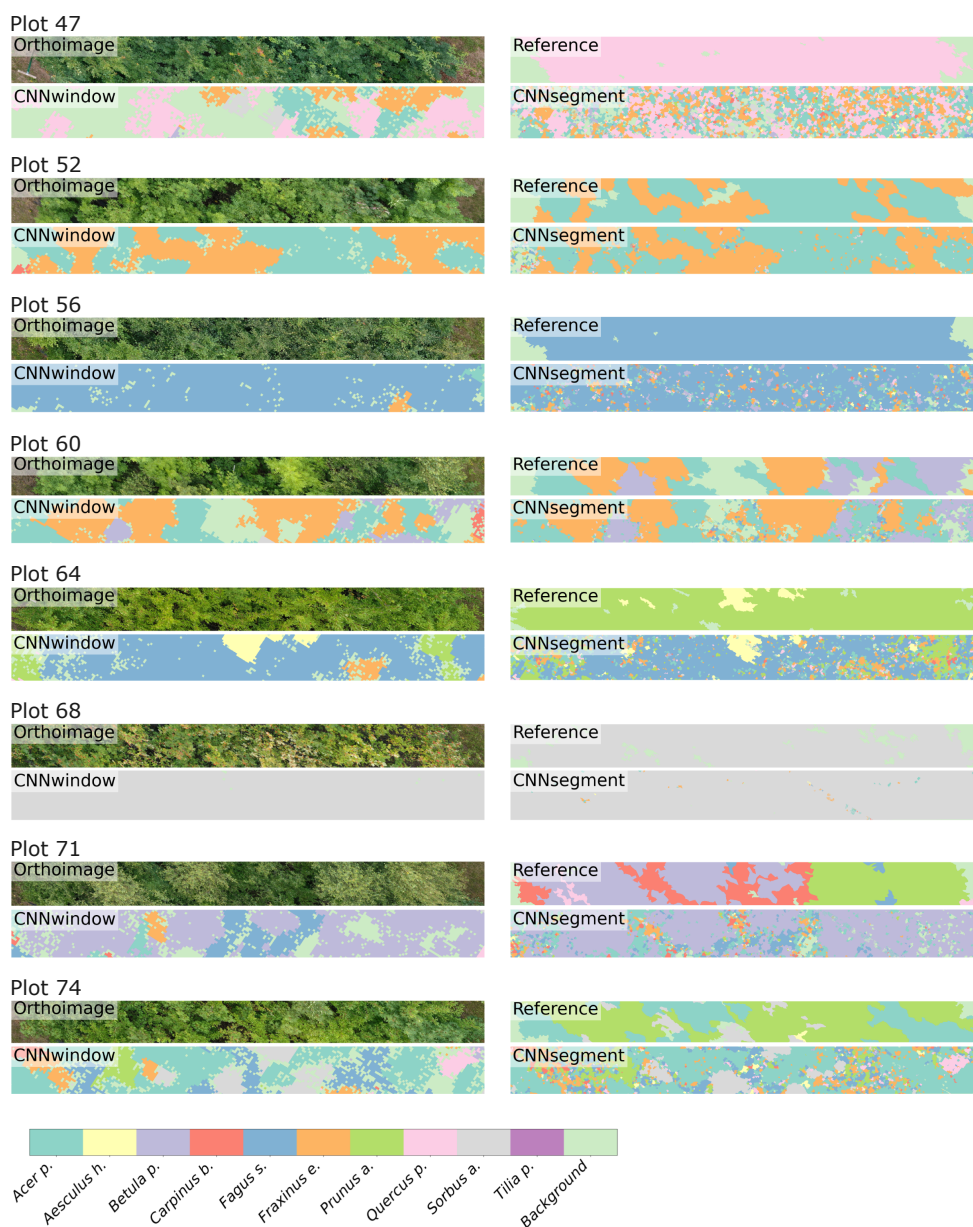


Figure A8: Transects of 2 m by 20 m of selected plots, including the orthoimage, the reference, CNN_{window} predictions, and CNN_{segment} predictions.

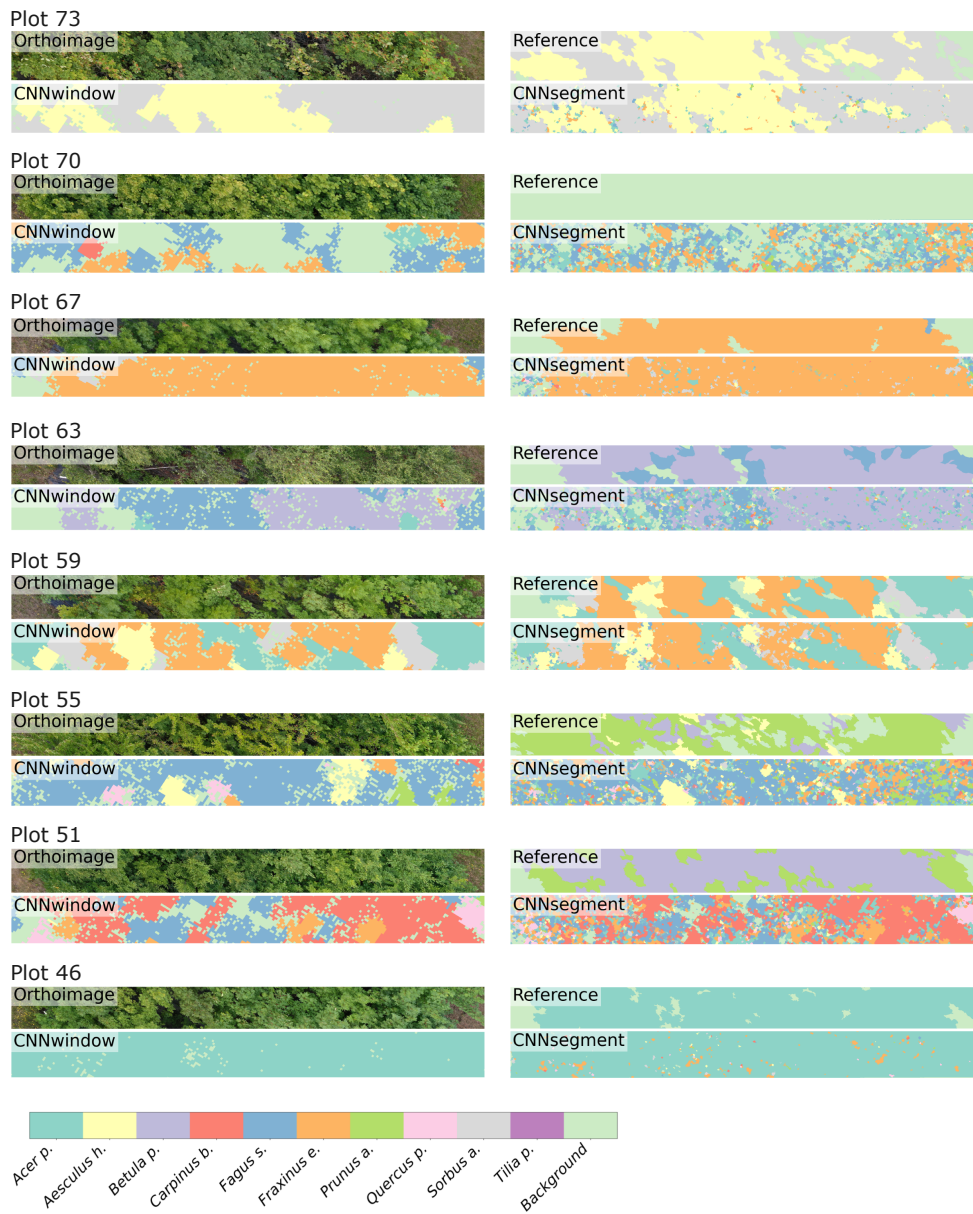


Figure A9: Transects of 2 m by 20 m of selected plots, including the orthoimage, the reference, CNN_{window} predictions, and $CNN_{segment}$ predictions.



650 **A1.2 Confusion Matrix**

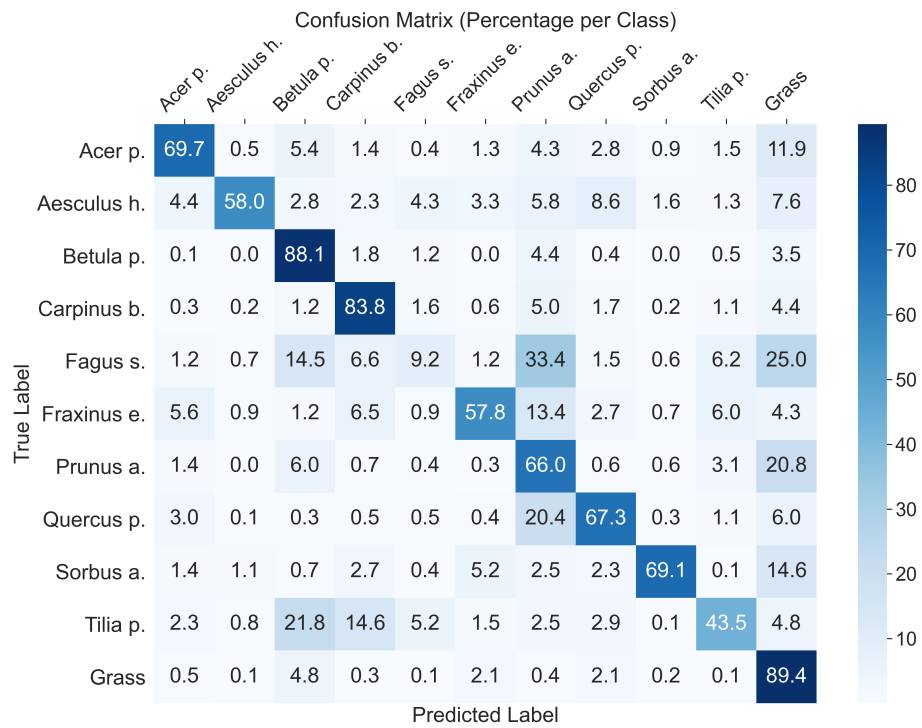


Figure A10: Normalized Confusion Matrix of the CNNsegment model applied to Ortho_{September}

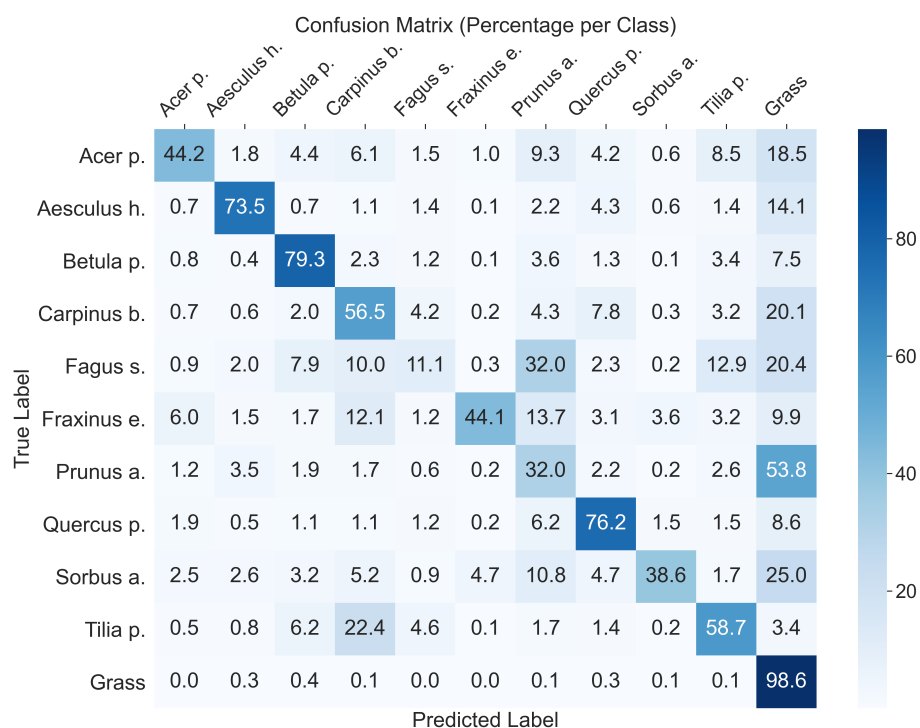


Figure A11: Normalized Confusion Matrix of the CNNsegment model applied to the Ortho_{September}

651 A1.3 Data pre-processing

652 To reduce the heterogeneity of crowd-sourced photographs and match them with the UAV
 653 perspective, we filtered the photographs based on their acquisition distance and plant leaf
 654 visibility. The model achieved an $R^2=0.7$ and $F1=0.8$ on independent test data for both
 655 variables. Using predicted acquisition distance and tree trunk presence information for each
 656 photograph, we tested different filtering thresholds and combinations prior to training the
 657 CNN_{window} model for plant species classification. The best result was achieved by filtering
 658 photographs with acquisition distances outside the range of 0.3 to 15 m and excluding pho-
 659 tographs that showed tree trunks, with a probability of being a trunk >0.5 .



660 **A1.4 Citizen science data availability**

Table A1: Number of downloaded photographs for selected tree species from the iNaturalist and Pl@ntNet datasets.

No.	Species	iNaturalist	Pl@ntNet
1	<i>Acer pseudoplatanus</i>	9999	3205
2	<i>Aesculus hippocastanum</i>	9998	1444
3	<i>Betula pendula</i>	9998	1308
4	<i>Carpinus betulus</i>	7165	2633
5	<i>Fagus sylvatica</i>	9981	3304
6	<i>Fraxinus excelsior</i>	7745	3130
7	<i>Prunus avium</i>	9999	3022
8	<i>Quercus petraea</i>	1491	221
9	<i>Sorbus aucuparia</i>	10000	2730
10	<i>Tilia platyphyllos</i>	582	1449

661 **A1.5 Segmentation model architecture**

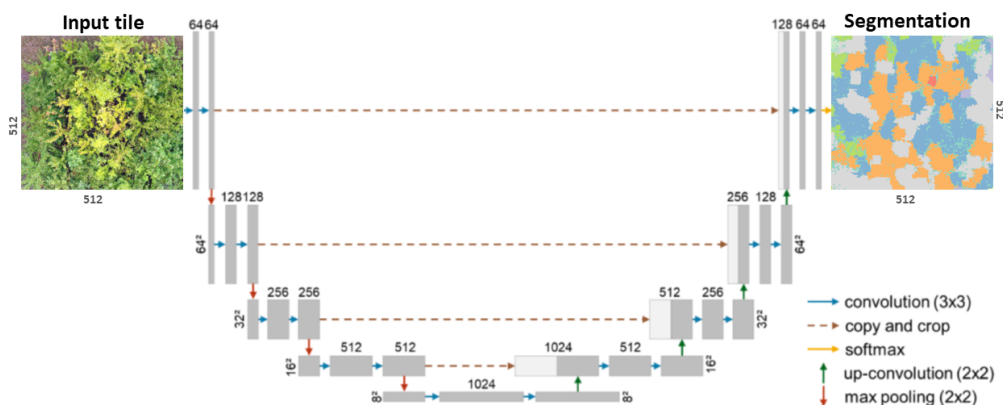
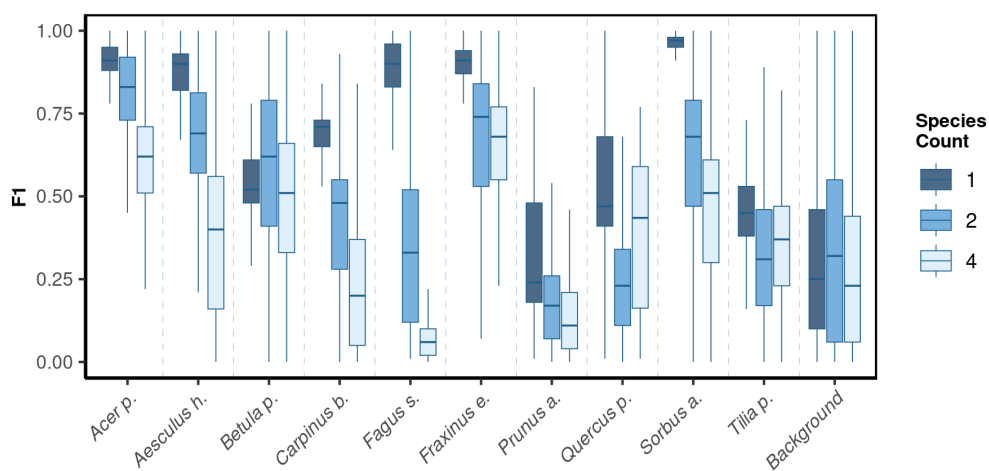


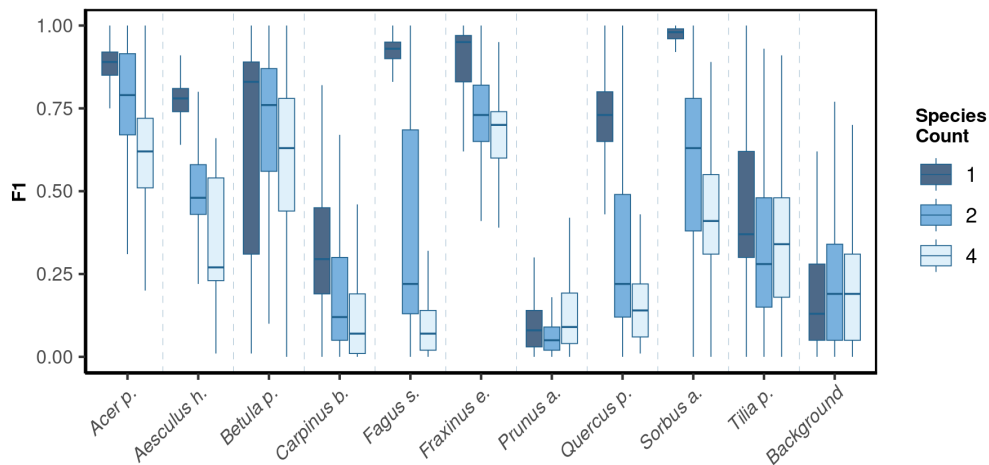
Figure A12: A modified version of the U-Net CNN-architecture for segmenting plant species from UAV orthoimages (Ronneberger et al., 2015).



662 **A1.6 CNN window species mixture box plot**



(a) Performance on Ortho_{July}: The model performance (F1) of the CNN_{window} model on Performance on Ortho_{July}.



(b) Performance on Ortho_{September}: The model performance (F1) of the CNN_{window} model on Performance on Ortho_{July}.

Figure A13: The model performance (F1) of the CNN_{segment} model across a gradient of canopy complexity in two orthoimages.