

1 **Title:** Habitat-Net: Segmentation of habitat images using deep learning

2 Jesse F. Abrams<sup>1\*†</sup>, Anand Vashishtha<sup>2\*</sup>, Seth T. Wong<sup>1</sup>, An Nguyen<sup>1</sup>, Azlan Mohamed<sup>1,3</sup>,  
3 Sebastian Wieser<sup>1</sup>, Arjan Kuijper<sup>2</sup>, Andreas Wilting<sup>1‡</sup>, Anirban Mukhopadhyay<sup>2‡</sup>

4 <sup>1</sup> Leibniz Institute for Zoo and Wildlife Research (IZW), 10315 Berlin, Germany

5 <sup>2</sup> Department of Computer Science, Technische Universität Darmstadt, 64283 Darmstadt,  
6 Germany

7 <sup>3</sup> WWF-Malaysia, Petaling Jaya, 46150, Selangor, Malaysia.

8 \* these authors contributed equally (co-first author)

9 ‡ these authors contributed equally (co-last author)

10 † corresponding author (jabrams23@gmail.com)

11 **ABSTRACT**

12 Understanding environmental factors that influence forest health, as well as the occurrence  
13 and abundance of wildlife, is a central topic in forestry and ecology. However, the manual  
14 processing of field habitat data is time-consuming and months are often needed to progress  
15 from data collection to data interpretation. Computer-assisted tools, such as deep-learning  
16 applications can significantly shortening the time to process the data while maintaining a  
17 high level of accuracy. Here, we propose Habitat-Net: a novel method based on  
18 Convolutional Neural Networks (CNN) to segment habitat images of tropical rainforests.  
19 Habitat-Net takes color images as input and after multiple layers of convolution and  
20 deconvolution, produces a binary segmentation of the input image. We worked on two  
21 different types of habitat datasets that are widely used in ecological studies to characterize  
22 the forest conditions: canopy closure and understory vegetation. We trained the model with  
23 800 canopy images and 700 understory images separately and then used 149 canopy and  
24 172 understory images to test the performance of Habitat-Net. We compared the  
25 performance of Habitat-Net with a simple threshold based method, a manual processing by a  
26 second researcher and a CNN approach called U-Net upon which Habitat-Net is based.  
27 Habitat-Net, U-Net and simple thresholding reduced total processing time to milliseconds per  
28 image, compared to 45 seconds per image for manual processing. However, the higher  
29 mean Dice coefficient of Habitat-Net (0.94 for canopy and 0.95 for understory) indicates that  
30 accuracy of Habitat-Net is higher than that of both the simple thresholding (0.64, 0.83) and  
31 U-Net (0.89, 0.94). Habitat-Net will be of great relevance for ecologists and foresters, who  
32 need to monitor changes in their forest structures. The automated workflow not only reduces  
33 the time, it also standardizes the analytical pipeline and, thus, reduces the degree of  
34 uncertainty that would be introduced by manual processing of images by different people  
35 (either over time or between study sites). Furthermore, it provides the opportunity to collect  
36 and process more images from the field, which might increase the accuracy of the method.  
37 Although datasets from other habitats might need an annotated dataset to first train the  
38 model, the overall time required to process habitat photos will be reduced, particularly for  
39 large projects.

40 **Index Terms:** Habitat Interpretation, Image Segmentation, Convolutional Neural Network,  
41 Deep Learning, canopy closure, understory vegetation density, forest

## 42 1. INTRODUCTION

43 Understanding of the intricacies of natural forest ecosystems is important to better  
44 manage and protect them. In both ecology and forestry there is a huge need for high quality  
45 information about the forest structure to help understand the spatio-temporal changes of  
46 forest habitat (Stojanova *et al.*, 2010). Detailed, large-scale knowledge about forest habitat  
47 and structure and how forests respond to anthropogenic disturbance will improve our ability  
48 to mitigate the effects of disturbances. Canopy closure and understory vegetation density are  
49 measurements that are commonly used in forest and land use research, monitoring,  
50 management and planning (Jennings, Brown, & Sheil, 1999).

51 Canopy closure is the proportion of sky hemisphere obscured by vegetation when  
52 viewed from a single point (Jennings, Brown, & Sheil, 1999). More rigorously defined,  
53 canopy closure is defined as the percent forest area occupied by the vertical projection of  
54 tree crowns (Paletto & Tosi, 2009). The canopy closure metric is an important part of forest  
55 inventories (Korhonen *et al.*, 2006; Chopping *et al.*, 2008), linked with canopy architecture,  
56 light regimes, solar radiation and leaf area index estimates in forest ecosystems. It is useful  
57 for wildlife habitat assessment and monitoring (Paletto & Tosi, 2009) and is often used as a  
58 multipurpose ecological indicator (Korhonen *et al.*, 2006). For forestry practitioners,  
59 measurements of the forest canopy serve as one of the chief indicators of the microhabitat  
60 within the forest (Jennings, Brown, & Sheil, 1999). The forest canopy affects plant growth  
61 and survival, hence determining the nature of the vegetation, and wildlife habitat.  
62 Characterization of the understory vegetation is of equal importance in forestry as it plays a  
63 central role in forest ecosystem structure and composition (Russell *et al.*, 2014). Understory  
64 vegetation provides key elements (and indicators) for biodiversity, nutrient cycling and forest  
65 fuel loads, and shape overstory tree structure and diversity (Halpern and Spies, 1995,  
66 Legare *et al.*, 2002; Gilliam, 2007; Russell *et al.*, 2014). Understory vegetation has become a  
67 fundamental component of forest site classifications (Bergès Gégout, & Franc, 2006) and the  
68 status of the understory composition and structure is a critical indicator of the condition of the  
69 forest (D'Amato, Orwig, Foster, 2009).

70 For ecologists and wildlife managers there is also a great need to understand factors  
71 influencing the occurrences and habitat preferences of species (Cristescu & Noyce, 2013),  
72 as the ecology of many species is poorly understood hindering their effective management.  
73 Therefore, accurate and quick habitat characterizations of the surveyed sites are needed  
74 (Zeng *et al.*, 2013). Canopy closure and understory vegetation are among the factors that

75 are known to influence species occurrence at a site (Vickers & Palmer, 2000;  
76 Brenes-Arguedas *et al.*, 2011).

77         Despite the importance of quantitative estimates of canopy closure and understory  
78 vegetation density, there are no efficient ways to obtain these values. A plethora of different  
79 techniques have been developed to quantify forest canopy (Jennings, Brown, & Sheil, 1999).  
80 Conventionally, information on canopy closure was collected using a spherical densiometer  
81 (Jennings, Brown, & Sheil, 1999). This technique is very labour intensive and requires  
82 researchers to spend a lot of time in the field. Currently, canopy closure information is often  
83 obtained through manual processing of color digital canopy photographs to binary images of  
84 vegetation and sky or through simple thresholding methods (Jonckheere *et al.*, 2005; Nobis,  
85 & Hunziker, 2005), which is less time consuming in the field but requires lots of manual  
86 processing at the computer. One of the traditional methods to characterize understory  
87 vegetation in the field is through the use of cover boards. Here an observer visually  
88 estimates the relative proportion of a board of known dimensions that is being obscured by  
89 vegetation from a given vantage point (Jones, 1968; Nudds, 1977). The subjectivity inherent  
90 to the visual estimation by an observer is a widely acknowledged limitation of both the  
91 spherical densiometer and cover board field methods (Limb *et al.*, 2007; Morrison, 2016).  
92 Recently, the workflow for the processing of both canopy and understory images has  
93 become digital due to advancements in digital photography and image processing of digital  
94 vegetation photographs (Marsden *et al.*, 2002; Jorgensen *et al.*, 2013).

95         Currently, due to logistical and analytical challenges, and time consuming manual  
96 processing of field data, several months, and sometimes up to years, are needed to progress  
97 through the stages of data collection, data processing, and interpretation. The processing of  
98 the habitat photograph datasets, which are usually large, is presently done using manual  
99 methods or simple thresholding methods. This is highly unsatisfactory and prevents timely  
100 evidence-based, effective action in both forestry and conservation applications. Fast,  
101 reliable, and automated computer-assisted tools are therefore needed to describe the habitat  
102 immediately after data collection. Using tools such as advanced machine learning, which  
103 have gained popularity when solving many data driven tasks in other fields, is one option to  
104 overcome the time demand associated with habitat interpretation. Some more advanced  
105 techniques for the processing of habitat images, some of which exploit machine learning, do  
106 exist, including traditional LiDAR-based and 3D image based techniques to segment canopy  
107 and understory images (Stojanova *et al.*, 2010; Tao *et al.*, 2015; Hamraz *et al.*, 2017). These  
108 methods however, require expensive equipment to collect LiDAR or 3D images and  
109 seldomly are these data available for forest inventories or ecological studies. Therefore most

110 projects still rely on either labour intensive field methods, such as the use of spherical  
111 densimeters, or on digital habitat photos, which later need manual processing. Thus, the  
112 processing of thousands of simple digital habitat photographs desperately requires advances  
113 in automated workflows.

114 Recently, computer vision has made an inroads to the ecological domain. There have  
115 been limited attempts to use machine learning methods for automated interpretation of forest  
116 habitat images, such as canopy closure photographs (Levner & Bulitko, 2004; Zhao *et al.*,  
117 2010; Erfanifard, Khodaei, & Shamsi, 2014; Ahmed *et al.*, 2015). Compared to these early  
118 studies, more advanced deep learning techniques have been developed, particularly in the  
119 field of medical imaging research. These methods have proved to be superior to earlier  
120 techniques (Litjeans *et al.*, 2017), significantly shortening the time and increasing the  
121 accuracy of the data interpretation and data processing. Specifically, Convolutional Neural  
122 Networks (CNNs), which are deep feedforward neural networks that are inspired by the  
123 visual cortex of the human eye and allow computers to ‘see’, have been shown to outperform  
124 many state-of-the-art methods in various visual computing tasks across different domains  
125 (Krizhevsky, Sutskever, & Hinton, 2012). CNNs are used to recognize images by processing  
126 the original image through multiple layers of feature-detecting “neurons”. Each layer is  
127 designed to detect a specific set of features such as lines, or edges. Increasing the number  
128 of layers (typical CNNs use anywhere from 5 to 25 layers or more) allows the CNN to detect  
129 more complex features enabling it to recognize the object in the image. Biomedical images  
130 generated by a wide range of medical procedures, such as MRI and CT scans, have  
131 complex textural patterns and are limited to small annotated datasets making it difficult to  
132 apply machine learning techniques and classic CNNs. However, the U-Net model helped  
133 overcome these challenges in biomedical image segmentation (Ronneberger, Fischer, &  
134 Brox, 2015). This is due to the network architecture that consists of a contracting path, in  
135 which the spatial information is reduced while feature information is increased, and an  
136 expansive path, which combines the feature and spatial information from the contracting  
137 path.

138 In practice, the application of CNNs to medical imaging is very similar to their  
139 application in the field of ecology. Similar to medical images, canopy and understory images  
140 contain a lot of color and textural information. In this work we focus on using computer vision  
141 methods to automatically extract the relevant information about canopy closure and  
142 understory vegetation density from in-situ digital photographs. We propose Habitat-Net, a  
143 novel deep learning method based on U-Net, to segment in-situ habitat images of forests.  
144 We extend the U-Net architecture to a new domain and improve on the performance of

145 U-Net by implementing Batch Normalization. Our proposed framework has been designed to  
146 feed color forest habitat (both canopy and understory) images as input to the network and,  
147 after multiple convolutions, generates a binary segmentation raster of the original image. The  
148 entire pipeline of the proposed method has been designed to work automatically without any  
149 user interaction. Our only assumption is that during model training the input images are  
150 accompanied by respective annotated images.

## 151 **2. METHODS**

### 152 ***2.1 Dataset Description***

153 We conducted standardized vegetation surveys in Sabah, Malaysian Borneo  
154 between 2014 and 2016. We collected a total of 949 canopy (128 x 128 pixels) and 872  
155 understory vegetation (256 x 160 pixels) photographs that are used in this study. All photos  
156 were taken using the built-in camera in GPS unit (Garmin® model 62sc). To collect our  
157 canopy dataset in the field, we established a 20 x 20 m grid around the center point of our  
158 survey station, which was located halfway between the two camera traps. The grid was  
159 positioned along the north-south, east-west axes. We took canopy photographs at the  
160 centerpoint and the NW, NE, SW, and SE corners of the plot. All canopy photos were taken  
161 at an angle of approximately 90 degrees (directly overhead). The understory dataset was  
162 collected by taking photos of a 1.5 x 1.0 m orange fly-sheet positioned 10 m in each cardinal  
163 direction while standing at the centerpoint of the survey grid. The vegetation covering the  
164 flysheet is used to estimate the understory vegetation density. The orange sheet used during  
165 data collection separates the understory areas from the background providing a means to  
166 segment the understory images. The photographs range in complexity from a completely  
167 uncovered orange sheet with no understory visible (the reference is an entirely white image)  
168 to images where the orange sheet is completely covered due to dense understory structure  
169 (reference segmentation image is an entirely black image).

### 170 ***2.2 Manual segmentation - Deriving canopy closure and vegetation density from field*** 171 ***photographs***

172 We used the free and open source image manipulation software Gimp to process the  
173 canopy and understory images. We set color thresholds and used the binary indexing  
174 feature in Gimp to convert the color canopy images into binary black and white images with  
175 black representing foliage and white sky (Fig. 1). The processing of the understory  
176 vegetation photos followed the same basic workflow, segmenting a binary image, but  
177 included an additional preprocessing step in which we cropped the image to the extent of the

178 orange flysheet that was being photographed from a 10 m distance. Similar to the canopy  
179 closure this workflow resulted in a binary (black and white) raster with black representing  
180 vegetation and white the orange flysheet (gaps in understory vegetation). Canopy closure  
181 and vegetation density was then calculated from the classified binary images by  
182 automatically counting black (vegetation) and white (non-vegetation) pixels using the  
183 following R script:

```
184 # load image as raster and convert to matrix
185 r <- raster(files.tmp[j])
186 r.matr <- as.matrix(r)

187 # set threshold value to split vegetation and sky/empty space: range = 0 - 256. 128 is medium grey (needed
188 because jpg creation from binary image introduces some artefacts which are then made binary again)
189 thresholdValue = 128,

190 # consider as vegetation all pixels with value < thresholdValue
191 fraction_vegetation <- as.numeric(mean(r.matr < thresholdValue))

192 # the rest is sky/empty space
193 fraction_sky <- 1 - fraction_vegetation
```

#### 194 *Second manual segmentation*

195 In order to compare the performance and consistency of manual segmentation by  
196 different researchers, a second independent researcher performed the manual segmentation  
197 on the test set of canopy and understory images.

### 198 **2.3 Simple thresholding**

199 Image thresholding is a simple, yet effective, image segmentation technique that  
200 partitions an image into a foreground and background. We ran a simple thresholding  
201 algorithm (scripts available in the Supplementary Material) to convert our color images to  
202 monochrome images. We first converted field images to grayscale and then used Otsu's  
203 method to automatically reduce the grayscale image to a binary image (Otsu, 1979). In short,  
204 in Otsu's method we assume that the image contains two classes of pixels following a  
205 bimodal histogram (foreground pixels and background pixels). We then calculate the  
206 threshold that minimizes the intra-class variance (the variance within the class), defined as a  
207 weighted sum of variances of the two classes by iterating through all possible threshold

208 values. We implemented this in Python version 3.7.1 using the *threshold\_otsu* function from  
209 the package *scikit\_image* version 0.14.1 (van der Walt *et al.*, 2014).

## 210 **2.4 Habitat-Net**

211 The Habitat-Net architecture (Fig. 2) is based on the U-Net (Ronneberger, Fischer, &  
212 Brox, 2015) convolutional network which provides pixel level localization by combining high  
213 resolution features with upsampling layer outputs. The U-Net model architecture has a large  
214 number of feature channels that allow the network to propagate context information to higher  
215 resolution layers. As a consequence, the expansive path is symmetric to the contracting path  
216 and yields a U-shaped architecture (Ronneberger, Fischer, & Brox, 2015).

217 The Habitat-Net consists of multiple  $3 \times 3$  convolutions followed by a non-linear  
218 activation using rectified linear unit (ReLU). The use of a small filter size helps to capture  
219 finer details of the image. To achieve a numerically stable training procedure, we incorporate  
220 a batch normalization (BN) layer (Ioffe & Szegedy, 2015) after every convolution layer as a  
221 novel design choice in Habitat-Net. The batch normalization layer reduces the internal  
222 covariate shift, which can boost segmentation performance and helps to make training more  
223 resilient to the parameter scale using mini-batches. Batch normalization further reduces  
224 overfitting to the minimum extent and replaces the need for a dropout layer in most cases  
225 (Srivastava *et al.*, 2019). Every batch normalization layer is followed by a  $2 \times 2$  max pooling  
226 operation, which operates independently on every depth slice of the input matrix and resizes  
227 it spatially using the max operation. We increased feature channels by an order of two at  
228 every downsampling step and halved the feature channels at each upsampling step. We  
229 used a zero-padding hyperparameter to control the spatial size of the output volumes  
230 (source code is provided in: <https://github.com/Kanvas89/Habitat-Net>).

231 Our final network includes a total of 33 convolutions including the final  $1 \times 1$  out  
232 convolution layer. The presence of many convolution layers helps to achieve better model  
233 accuracy (Szegedy *et al.*, 2015). The only trade-off of this increased accuracy is that the  
234 network requires more time and resources to converge. We use a stochastic gradient  
235 descent (SGD) optimizer with time-based decay, which handles limited datasets well. SGD  
236 performs better in vision-based machine learning tasks and generalizes quicker than other  
237 adaptive methods such as Adam, RMSProp, and AdaGrad (Wilson *et al.*, 2017). SGD with  
238 momentum helps the parameter vector to build up velocity in any direction with a constant  
239 gradient descent so as to prevent oscillations, which leads to faster-convergence (Qian,  
240 1999; Sutskever *et al.*, 2013). The final output layer returns a one-channel grayscale  
241 prediction image of the input habitat (canopy and understory) image. The binary raster is



242 then run through the same R function that is used for the manual processing. All the  
243 experiments have been run using Keras 2.2.2 with TensorFlow 1.10.0 using python 3.5 on a  
244 system with a Nvidia 1080 Ti GPU.

#### 245 *2.4.1 Network training and testing*

246 Here, we test our method using the canopy closure and understory density datasets.  
247 The training datasets consist of a pair of images (either canopy or understory), the image to  
248 be segmented and a manual segmentation raster drawn by an expert to be used as a  
249 reference to train the model (Fig. 1). Deep neural networks typically perform better with more  
250 training data. Models trained on small datasets do not generalize well and suffer from  
251 overfitting (Perez & Wang, 2017). When a limited number of images with complex textural  
252 and color patterns are available to train the model, it is imperative to exploit data  
253 augmentation to increase the total number of training images. However, the non-linear  
254 transformations used for augmenting cell images in U-Net (Ronneberger, Fischer, & Brox,  
255 2015) can not be incorporated directly in Habitat-net due to domain specific properties. We  
256 address this issue by artificially inflating the number of training images through rotations and  
257 reflection of the image (Supplementary Figs. S1 & S2).

258 Of 949 color canopy images, the training consists of 800 image pairs (image and  
259 respective reference segmentation raster). After applying the data augmentation technique  
260 our training dataset has a total of 4000 image pairs. We use 3400 images for the training  
261 dataset and the remaining 600 for validation during training. Of the 872 understory  
262 vegetation photos 700 images (after augmentation 3500 images) were used in the training  
263 dataset. Similar to the manual processing, it was necessary to perform the preprocessing  
264 step of cropping the understory images to the extent of the orange flysheet as this could not  
265 be automated. The code and the trained weights for Habitat-Net can be accessed at  
266 <https://github.com/Kanvas89/Habitat-Net>. As the network converged faster and variance was  
267 lower when batch normalization was used after each convolution layer (Supplementary Fig.  
268 S3) we applied batch normalization in all results presented below.

269 After training Habitat-Net, we tested its performance on the remaining 149 canopy  
270 images and 172 understory images (15% of the total dataset). To evaluate the performance  
271 of Habitat-Net we used overlap ratio measures (Jaccard 1907; Dice 1945; Sørensen, 1948),  
272 which quantify the degree of similarity between two objects. In our case they are an indicator  
273 of the overlap between our manually segmented reference images and those generated by  
274 the Habitat-Net model. We report both the Dice coefficient and Jaccard index as  
275 segmentation quality metrics to evaluate the performance of Habitat-Net. This is because

276 although the Dice coefficient and Jaccard index are similar, the Jaccard index is numerically  
277 more sensitive to mismatch when there is reasonably strong overlap. As the Jaccard index  
278 penalizes single instances of bad classification more than the Dice coefficient, results of the  
279 Dice coefficient typically "look nicer" because they are higher for the same pair of  
280 segmentations and thus the Dice coefficient index is currently more popular than the Jaccard  
281 index.

### 282 **3. RESULTS**

283 Based on researcher experience, the manual processing of a canopy closure image  
284 requires about 45 seconds per image, while the processing of an understory image requires  
285 about 65 seconds per image (this largely varies depending on the level of experience of the  
286 individual researcher and quality of the images). Habitat-Net reduces total processing time to  
287 around 15 milliseconds per image for both canopy and understory images. Similar simple  
288 thresholding significantly reduces the time for 1 image to less than 1 second. For a typical  
289 dataset of around 400 images, both Habitat-Net and simple thresholding reduce total  
290 processing time to seconds compared to the manual processing for which 5 hours (canopy  
291 images) or even 7.5 hours (understory images) were needed (Tables 1 and 2). Visual  
292 inspection of the segmentation results (Fig. 3) from the three automated methods indicate  
293 that Habitat-Net and U-Net (both machine learning methods) outperform the simple  
294 thresholding method for the canopy images. However, the improved performance for the  
295 understory images is not as obvious in many cases. However the quantitative assessment of  
296 the performance of the different methods using the Dice and Jaccard similarity scores reveal  
297 the greatest accuracy of Habitat-Net (Fig. 4, Tables 1 & 2). Particular for the canopy images  
298 the similarity scores of the Habitat-Net were with 0.94 Dice and 0.88 (Jaccard) much higher  
299 than for the other methods, including the U-Net upon which Habitat-Net is based (Table 1).  
300 The differences for the understory images between the different methods were less strong,  
301 as all methods had higher similarity scores. However again Habitat-Net outperformed the  
302 other methods with a Dice score of 0.95 and a Jaccard index of 0.92 (Table 2). Although  
303 Habitat-Net generally outperformed other segmentation methods, there were a few extreme  
304 outliers produced. We visually inspected the color photographs of the images with outlying  
305 similarity scores for Habitat-Net (Fig. 4). For canopy images, the images with the poorest  
306 similarity scores are "speckled" images that contain many small (sometimes single pixel)  
307 openings in the canopy vegetation. However, for the canopy images the minimal similarity  
308 score is 0.75, indicating that even when Habitat-Net performs poorly, the resulting  
309 segmentation is still acceptable. There are, however, some more extreme outliers present for

310 the understory images. These images are very dark, blurry, and either all or almost all of the  
311 orange flysheet is covered, leaving only small openings in the vegetation. Although we  
312 inspected the problems in images which were outliers in the Habitat-Net analysis similar  
313 outliers, very likely in the same images were found in all methods, even in the manual  
314 processing of a second researcher.

#### 315 **4. DISCUSSION**

316 Deep convolutional networks have been shown to outperform many other methods in  
317 various visual computing tasks and domains (Krizhevsky, Sutskever, & Hinton, 2012). CNNs  
318 have, however, seen little application in the ecological domain. With Habitat-Net we present  
319 an automated pipeline to process hundreds of color vegetation photographs (both canopy  
320 and understory) in a standardized, efficient and reproducible way. Our approach saves a  
321 huge amount of human labor and helps overcome the time demand associated with habitat  
322 interpretation. Habitat-Net has two advantages over manual processing, segmentation is: (1)  
323 significantly faster, and (2) more consistent. Habitat-Net produces binary segmentations with  
324 a higher similarity to the manual reference segmentation than do the approaches using  
325 simple thresholding or U-Net. For canopy images, the implementation of Habitat-Net led to a  
326 significant increase in accuracy and consistency of the image segmentation. For the  
327 understory images all methods produced a high similarity score, with Habitat-Net performing  
328 best and edging out U-Net. Habitat-Net performed well to segment images with a wide range  
329 of lighting conditions and sharpness, and the proposed method performed remarkably well  
330 even in situations where very little sky or orange flysheet was visible in a photograph (Fig. 3).  
331 The inclusion of a batch normalization layer after every convolution layer in Habitat-Net  
332 avoided internal covariate shift and enabled faster learning rates. This proved to stabilize the  
333 training (Fig. S3), boost the performance of Habitat-Net and improve the accuracy of  
334 multi-channel segmentation. Habitat-Net also outperforms previous machine learning-based  
335 methods for quantifying canopy closure. Previous research applying the ADaptive Object  
336 REcognition (ADORE; Draper, Bins, & Baek, 2000) system to canopy closure segmentation  
337 tasks produced a mean pixel-level similarity score (Dice coefficient) of  $0.54 \pm 0.14$  (Levner &  
338 Bilitko, 2004). In contrast to canopy closure, so far there are no other automated pipelines  
339 available for the analysis of understory photographs. Therefore, Habitat-Net is the first tool  
340 that allows foresters and ecologist to describe quantitatively both the horizontal and vertical  
341 forest structures.

342 In this study we only quantify understory vegetation density and we did not assess  
343 vegetation complexity. However, the quantification of vegetation complexity, which provides

344 further insight into forest structure, and thus into the ecosystem function of the understory  
345 (Halpern and Spies, 1995, Legare *et al.*, 2002; Gilliam, 2007; Russell *et al.*, 2014) would also  
346 be possible using the binary raster produced by Habitat-Net. In this case, researchers would  
347 need to take many field photographs with the contrasting flysheet at different distances. As  
348 the processing of these hundreds of photographs is automated with Habitat-Net these photo  
349 series could be used to reconstruct the vegetation complexity, without the need of expensive  
350 ground-based or airborne LiDAR scans.

351 For ecologists, the canopy and understory habitat are important indicators of forest  
352 disturbance and, for some wildlife species, tracking these disturbances might be a warning  
353 signal of potential population declines. Furthermore, the habitat information can be combined  
354 with spatial statistics, such as species distribution models (SDMs) to assess species  
355 occurrence or abundance data (Niedballa *et al.*, 2015). This allows researchers to determine  
356 habitat associations of little known species and the knowledge gained ecological about the  
357 species can be used for more effective conservation efforts. Therefore, Habitat-Net has the  
358 potential to be an efficient and effective tool for both foresters and wildlife ecologists.

359 Although our network is designed with small datasets in mind, deep learning works  
360 best with large datasets (Goodfellow, Bengio, & Courville, 2016; Norouzzadeh *et al.*, 2018).  
361 Even with a relatively small training dataset available, Habitat-Net performed with a high  
362 level of accuracy. Norouzzadeh *et al.* (2018) point out that the accuracy of deep learning  
363 methods further improves as more labeled data are provided during training. Thus, as more  
364 datasets become available the performance of the network may improve by building on  
365 knowledge from multiple datasets in a process known as transfer learning (Norouzzadeh *et al.*  
366 *et al.*, 2018).

367 A major limitation of manually processed images is the lack of standardization. The  
368 manual processing of images by humans introduces observer bias. For example, within one  
369 project a few people may do the manual processing, or in a long term monitoring program  
370 the observer (for example forestry staff) changes with time. Both of these scenarios would  
371 introduce bias and the subjectiveness inherent in the manual processing of images, making  
372 comparisons between the photographs or years difficult. Our study showed that the similarity  
373 scores between two different researchers doing the manual processing were lower than  
374 using Habitat-Net. One of the biggest strengths of Habitat-Net is, therefore, the ability to  
375 eliminate this user produced bias and standardize results which allows for inter-study site  
376 comparisons and long term monitoring.

377 Automated workflows, by default, accept biased inputs and can, therefore, generate  
378 undesired results. During manual processing the observer is able to sort out photos which

379 are out of focus or of poor quality. These photos are, however, included in the automated  
380 workflow, which may lead to inaccurate estimations. Such poor quality photos could  
381 potentially be automatically removed, but this would require a high level of standardization in  
382 how the photos are taken in the field. For all photos we used the GPS unit's built-in camera,  
383 which is a very basic camera that allows little control over camera settings. To increase  
384 standardization to a level that could potentially allow for automated removal of poor quality  
385 photos all images would have to be taken with the same camera, same settings and at the  
386 same camera angle, framing and distance. Although this would require carrying an additional  
387 higher quality camera into the field, as well as more time spent to set up equipment, we are  
388 certain that such higher quality photographs will increase the performance of Habitat-Net  
389 and, thus, the accuracy of the analysis. Currently it is not possible to automatically crop the  
390 understory vegetation photos and, thus, each image still required a minimal amount of  
391 manual processing. A satisfactory solution to do this for our images could not be found as  
392 often most of the orange flysheet was covered by vegetation making the automated  
393 recognition of the flysheet impossible. Common methods used to automatically crop images  
394 cannot be applied to our understory vegetation photographs due to the lack of common  
395 features that are always present (such as a frame). Other machine learning methods for  
396 cropping focus on "salient" image regions. The basic idea is to use information learned by  
397 the CNN about where human viewers fix their gaze to center a crop around the most  
398 interesting region (Rahman *et al.*, 2018). However, these methods are also not applicable for  
399 our understory photographs since the photographs are cluttered with no one point of interest.  
400 Such limitations could be overcome through innovative standardized practices in the field.  
401 The best solution would be to always have the flysheet centered and then perform a batch  
402 crop on all images to the specified area. This could be implemented in the field by placing a  
403 "stencil" or guide over the camera lens that is used to frame the flysheet in the center of the  
404 image. Then a technique called Object Localization and Detection could be implemented to  
405 detect the bounding box or frame (Sermanet *et al.*, 2013).

406 Habitat-Net provides a fast, accurate and standardized method to analyse canopy  
407 closure and understory vegetation photographs. With some optimisations during the field  
408 data collection, such as using a higher quality camera, placing a stencil over the camera lens  
409 the accuracy of Habitat-Net could even be improved and the time consuming manual  
410 cropping would not be necessary any more. With Habitat-Net, we can overcome the time lag  
411 from data collection to data processing, which often hinders timely management decisions  
412 and thus assist more sustainable forest management and conservation. Studies in other  
413 habitats might need a preliminary annotated dataset to train the model, but the overall time

414 required to process habitat photos will be reduced and the accuracy will be increased,  
415 particularly for large projects. We hope that other users add additional datasets and their  
416 modifications of the codes to github to expand the applications and focus of Habitat-Net  
417 further. Through this a large repository of habitat images can be built, which would in turn  
418 benefit both the ecology and machine learning communities.

### 419 **Acknowledgements**

420 We thank the Sabah Biodiversity Center and the Sabah Forestry Department, especially  
421 Johnny Kissing and Peter Lagan for support and involvement in this project. Many thanks go  
422 to the field team for their hard work to take all the habitat photographs. We would also like to  
423 thank Srijita Guha for helping with the manual cropping of understory images. This project  
424 received financial support from the German Federal Ministry of Education and Research  
425 (BMBF FKZ: 01LN1301A), Point Defiance Zoo and Aquarium through Dr. Holly Reed  
426 Conservation Fund and San Francisco Zoo.

### 427 **Conflict of interest declarations**

428 The other authors declare no conflicts of interests.

### 429 **Author contributions**

430 AW and AM designed the study; AV, AK, and AM wrote, trained, and test the Habitat-Net  
431 model. STW and AM collected the field data. STW, AN, AM, and SW conducted the manual  
432 image segmentation. JFA analyzed the results. JFA, AW, and AM lead writing the  
433 manuscript. All authors contributed to drafts and gave final approval for publication.

### 434 **References**

- 435 Ahmed, O. S., Franklin, S. E., Wulder, M. A., & White, J. C. (2015). Characterizing stand-level forest canopy  
436 cover and height using landsat time series, samples of airborne LiDAR, and the random forest  
437 algorithm. *ISPRS Journal of Photogrammetry and Remote Sensing*, 101, 89-101.
- 438 Bergès, L., Gégout, J. C., & Franc, A. (2006). Can understory vegetation accurately predict site index? A  
439 comparative study using floristic and abiotic indices in sessile oak (*Quercus petraea* Liebl.) stands in  
440 northern France. *Annals of forest science*, 63(1), 31-42.
- 441 Brenes-Arguedas, T., Roddy, A.B., Coley, P.D., & Kursar, T.A. (2011). Do differences in understory light contribute  
442 to species distributions along a tropical rainfall gradient? *Oecologia* 166-166(2), 443-456.
- 443 Chopping, M., Moisen, G. G., Su, L., Laliberte, A., Rango, A., Martonchik, J. V., & Peters, D. P. (2008). Large  
444 area mapping of southwestern forest crown cover, canopy height, and biomass using the NASA  
445 Multiangle Imaging Spectro-Radiometer. *Remote Sensing of Environment*, 112(5), 2051-2063.
- 446 Cristescu, B. & Boyce, M.S. (2013). Focusing ecological research for conservation. *Ambio* 42(7), 805- 815.

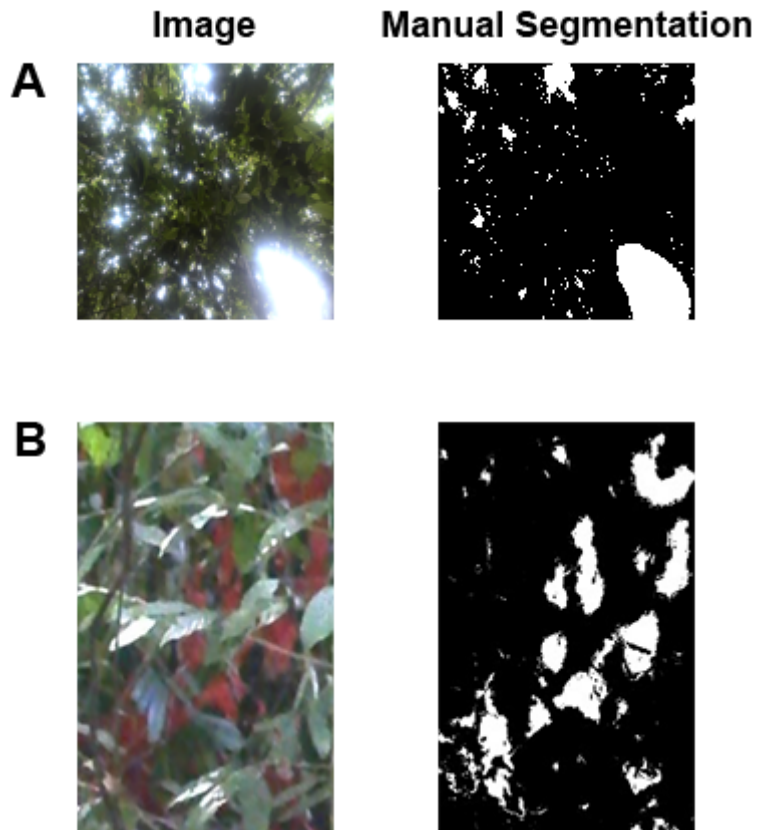
- 447 D'Amato, A. W., Orwig, D. A., & Foster, D. R. (2009). Understory vegetation in old-growth and second-growth  
448 *Tsuga canadensis* forests in western Massachusetts. *Forest Ecology and Management*, 257(3),  
449 1043-1052.
- 450 Dice, L. R. (1945). Measures of the amount of ecologic association between species. *Ecology*, 26(3), 297-302.
- 451 Draper, B. A., Bins, J., & Baek, K. (1999, January). ADORE: adaptive object recognition. In *International*  
452 *Conference on Computer Vision Systems* (pp. 522-537). Springer, Berlin, Heidelberg.
- 453 Erfanfard, Y., Khodaei, Z., & Shamsi, R. F. (2014). A robust approach to generate canopy cover maps using  
454 UltraCam-D derived orthoimagery classified by support vector machines in Zagros woodlands, West  
455 Iran. *European Journal of Remote Sensing*, 47(1), 773-792.
- 456 Gilliam, F. S., & Roberts, M. R. (2003). Interactions between the herbaceous layer and overstory canopy of  
457 eastern forests. The herbaceous layer in forests of eastern North America. Oxford University Press,  
458 Oxford, UK, 198-223.
- 459 Gilliam, F. S. (2007). The ecological significance of the herbaceous layer in temperate forest ecosystems. *AIBS*  
460 *Bulletin*, 57(10), 845-858.
- 461 Goodfellow, I., Bengio, Y., Courville, A., & Bengio, Y. (2016). Deep learning (Vol. 1). Cambridge: MIT press.
- 462 Halpern, C. B., & Spies, T. A. (1995). Plant species diversity in natural and managed forests of the Pacific  
463 Northwest. *Ecological Applications*, 5(4), 913-934.
- 464 Hamraz, H., Contreras, M. A., & Zhang, J. (2017). Vertical stratification of forest canopy for segmentation of  
465 understory trees within small-footprint airborne LiDAR point clouds. *ISPRS Journal of Photogrammetry*  
466 *and Remote Sensing*, 130, 385-392.
- 467 Ioffe, S. & Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal  
468 covariate shift. *arXiv preprint arXiv:1502.03167*.
- 469 Jaccard, P. (1907). La distribution de la flore dans la zone alpine. *Revue generale des Sciences pures et*  
470 *appliquees*, 18(23), 961-967
- 471 Jennings, S. B., Brown, N. D., & Sheil, D. (1999). Assessing forest canopies and understorey illumination: canopy  
472 closure, canopy cover and other measures. *Forestry*, 72(1), 59-74.
- 473 Jonckheere, I., Nackaerts, K., Muys, B., & Coppin, P. (2005). Assessment of automatic gap fraction estimation of  
474 forests from digital hemispherical photography. *Agricultural and Forest Meteorology*, 132(1-2), 96-114.
- 475 Jones, R. (1968). Productivity studies on heath vegetation in southern Australia the use of fertilizer in studies of  
476 production processes. *Folia Geobotanica et Phytotaxonomica*, 3(4), 355-362.
- 477 Jorgensen, C.F., Stutzman, R.J., Anderson, L.C., Decker, S.E., Powell, L.A., Schacht, W.H., & Fontaine, J.J.  
478 (2013). Choosing a DIVA: a comparison of emerging digital imagery vegetation analysis techniques.  
479 *Applied Vegetation Science*, 16(4), 552-560.
- 480 Korhonen, L., Korhonen, K. T., Rautiainen, M., & Stenberg, P. (2006). Estimation of forest canopy cover: a  
481 comparison of field measurement techniques. *Silva Fennica*, 40(4), 577-588.
- 482 Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep convolutional neural  
483 networks. In *Advances in neural information processing systems* (pp. 1097-1105).
- 484 Légaré, S., Bergeron, Y., & Paré, D. (2002). Influence of forest composition on understory cover in boreal  
485 mixedwood forests of western Quebec. *Silva Fennica*, 36(1), 353-366.
- 486 Levner, I., & Bulitko, V. (2004, July). Machine learning for adaptive image interpretation. In AAAI (pp. 870-876).

- 487 Limb, R.F., Hickman, K.R., Engle, D.M., Norland, J.E., & Fuhlendorf, S.D. (2007). Digital photography: reduced  
488 investigator variation in visual obstruction measurements for southern tallgrass prairie. *Rangeland*  
489 *Ecology & Management*, 60(5), 548-552.
- 490 Litjens, G., Kooi, T., Bejnordi, B. E., Setio, A. A. A., Ciompi, F., Ghafoorian, M., ... & Sánchez, C. I. (2017). A  
491 survey on deep learning in medical image analysis. *Medical image analysis*, 42, 60-88.
- 492 Marsden, S. J., Fielding, A. H., Mead, C., & Hussin, M. Z. (2002). A technique for measuring the density and  
493 complexity of understorey vegetation in tropical forests. *Forest Ecology and Management*, 165(1-3),  
494 117-123.
- 495 Morrison, L.W. (2016). Observer error in vegetation surveys: a review. *Journal of Plant Ecology*, 9, 367-379.
- 496 Niedballa, J., Sollmann, R., bin Mohamed, A., Bender, J., & Wilting, A. (2015). Defining habitat covariates in  
497 camera-trap based occupancy studies. *Scientific Reports*, 5, 17041.
- 498 Nobis, M., & Hunziker, U. (2005). Automatic thresholding for hemispherical canopy-photographs based on edge  
499 detection. *Agricultural and Forest Meteorology*, 128(3-4), 243-250.
- 500 Norouzzadeh, M. S., Nguyen, A., Kosmala, M., Swanson, A., Palmer, M. S., Packer, C., & Clune, J. (2018).  
501 Automatically identifying, counting, and describing wild animals in camera-trap images with deep  
502 learning. *Proceedings of the National Academy of Sciences*, 201719367.
- 503 Nudds, T. D. (1977). Quantifying the vegetative structure of wildlife cover. *Wildlife Society Bulletin*, 113-117.
- 504 Otsu, N. (1979). A threshold selection method from gray-level histograms. *IEEE transactions on systems, man,*  
505 *and cybernetics*, 9(1), 62-66.
- 506 Paletto, A., & Tosi, V. (2009). Forest canopy cover and canopy closure: comparison of assessment techniques.  
507 *European Journal of Forest Research*, 128(3), 265-272.
- 508 Perez, L., & Wang, J. (2017). The effectiveness of data augmentation in image classification using deep learning.  
509 arXiv preprint arXiv:1712.04621.
- 510 Qian, N. (1999). On the momentum term in gradient descent learning algorithms. *Neural Networks*, 12(1),  
511 145-151.
- 512 Rahman, Z., Pu, Y. F., Aamir, M., & Ullah, F. (2018). A framework for fast automatic image cropping based on  
513 deep saliency map detection and gaussian filter. *International Journal of Computers and Applications*,  
514 1-11.
- 515 Ronneberger, O., Fischer, P., & Brox, T. (2015). U-net: Convolutional networks for biomedical image  
516 segmentation. In *International Conference on Medical image computing and computer-assisted*  
517 *intervention* (pp. 234-241). Springer, Cham.
- 518 Russell, M. B., D'Amato, A. W., Schulz, B. K., Woodall, C. W., Domke, G. M., & Bradford, J. B. (2014).  
519 Quantifying understorey vegetation in the US Lake States: a proposed framework to inform regional  
520 forest carbon stocks. *Forestry*, 87(5), 629-638.
- 521 Sermanet, P., Eigen, D., Zhang, X., Mathieu, M., Fergus, R., & LeCun, Y. (2013). Overfeat: Integrated recognition,  
522 localization and detection using convolutional networks. arXiv preprint arXiv:1312.6229.
- 523 Sørensen, T. (1948). A method of establishing groups of equal amplitude in plant sociology based on similarity of  
524 species and its application to analyses of the vegetation on Danish commons. *Biol Skr*, 5, 1-34.
- 525 Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2014). Dropout: a simple way to  
526 prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1), 1929-1958.
- 527 Stojanova, D., Panov, P., Gjorgjioski, V., Kobler, A., & Džeroski, S. (2010). Estimating vegetation height and  
528 canopy cover from remotely sensed data with machine learning. *Ecological Informatics*, 5(4), 256-266.

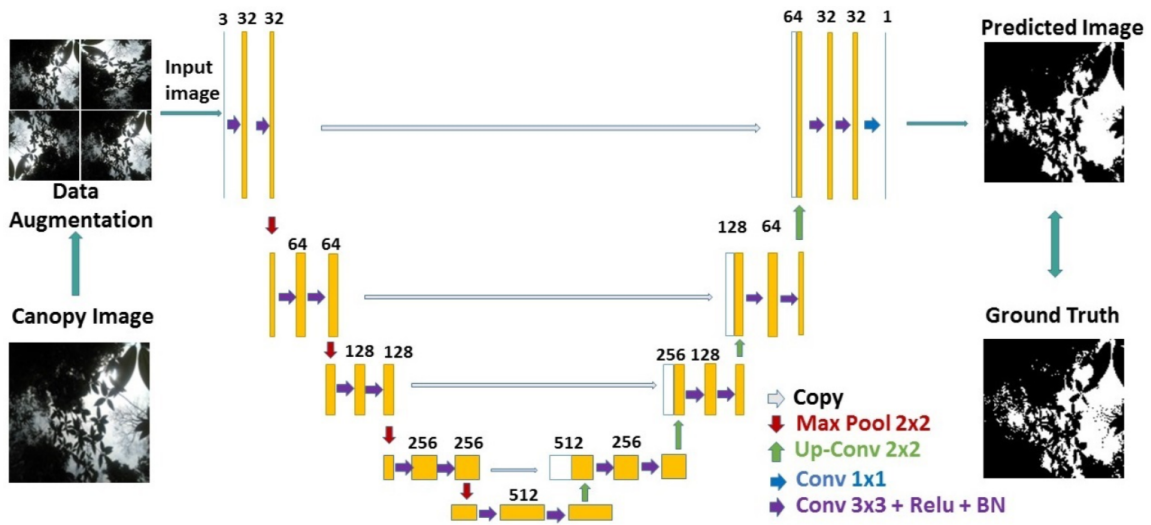


- 529 Sutskever, I., Martens, J., Dahl, G., & Hinton, G. (2013, February). On the importance of initialization and  
530 momentum in deep learning. In *International conference on machine learning* (pp. 1139-1147).
- 531 Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., & Rabinovich, A.  
532 (2015). Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and*  
533 *pattern recognition* (pp. 1-9).
- 534 Tao, S., Wu, F., Guo, Q., Wang, Y., Li, W., Xue, B., Hu, X., Li, P., Tian, D., Li, C., & Yao, H. (2015). Segmenting  
535 tree crowns from terrestrial and mobile LiDAR data by exploring ecological theories. *ISPRS Journal of*  
536 *Photogrammetry and Remote Sensing*, 110, 66-76.
- 537 Van der Walt, S., Schönberger, J. L., Nunez-Iglesias, J., Boulogne, F., Warner, J. D., Yager, N., Gouillart, E., &  
538 Yu, T. (2014). scikit-image: image processing in Python. *PeerJ*, 2, e453.
- 539 Vickers, A.D. & Palmer, S.C.F. (2000). The influence of canopy cover and other factors upon the regeneration of  
540 Scots pine and its associated ground flora within Glen Tanar National Nature Reserve. *Forestry*, 73(1):  
541 37-49.
- 542 Wilson, A. C., Roelofs, R., Stern, M., Srebro, N., & Recht, B. (2017). The marginal value of adaptive gradient  
543 methods in machine learning. In *Advances in Neural Information Processing Systems* (pp. 4148-4158).
- 544 Zeng, Y., Xu, Y., Wang, Y., & Zhou, C. (2013). Habitat Association and Conservation Implications of Endangered  
545 Francois' Langur (*Trachypithecus francoisi*). *PLoS ONE* 8(10): e75661.
- 546 Zhao, K., Popescu, S., Meng, X., Pang, Y., & Agca, M. (2011). Characterizing forest canopy structure with lidar  
547 composite metrics and machine learning. *Remote Sensing of Environment*, 115(8), 1978-1996.

548 **Figures**

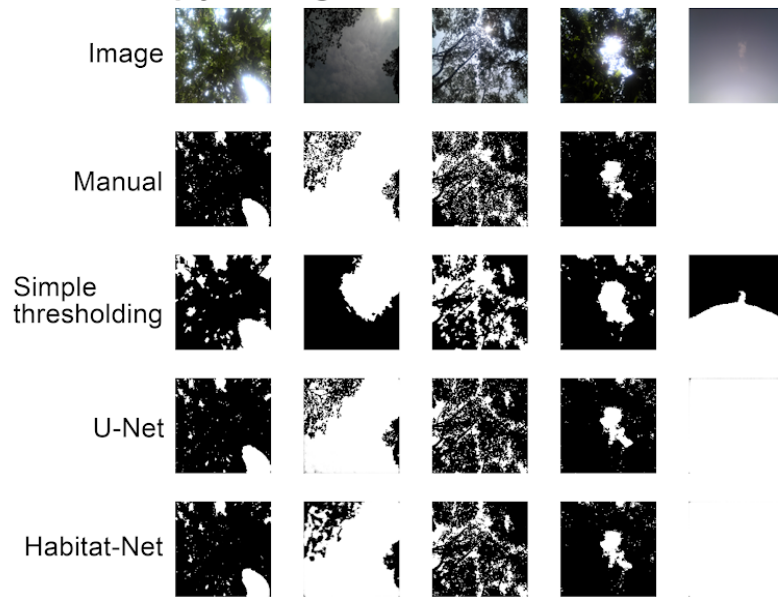


549 **Figure 1:** The first column contains the digital color images from photographs of (A) canopy  
550 and (B) understory taken in the field. The second column shows the manually segmented  
551 binary images of the same (A) canopy and (B) understory photographs.

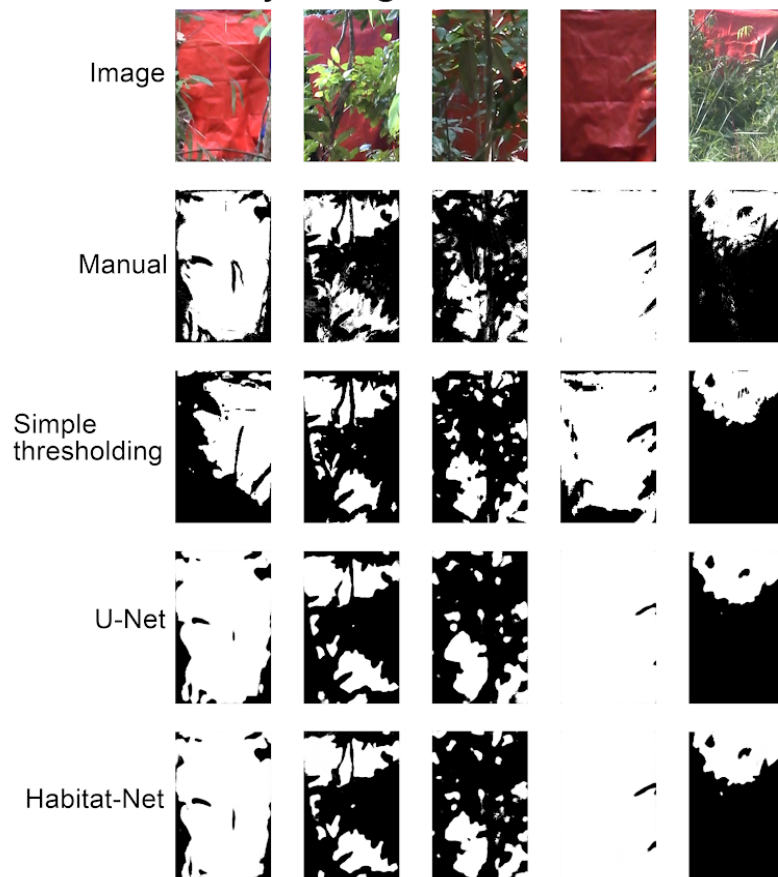


552 **Figure 2:** Habitat-Net architecture based on U-Net (Ronneberger, Fischer, & Brox, 2015).

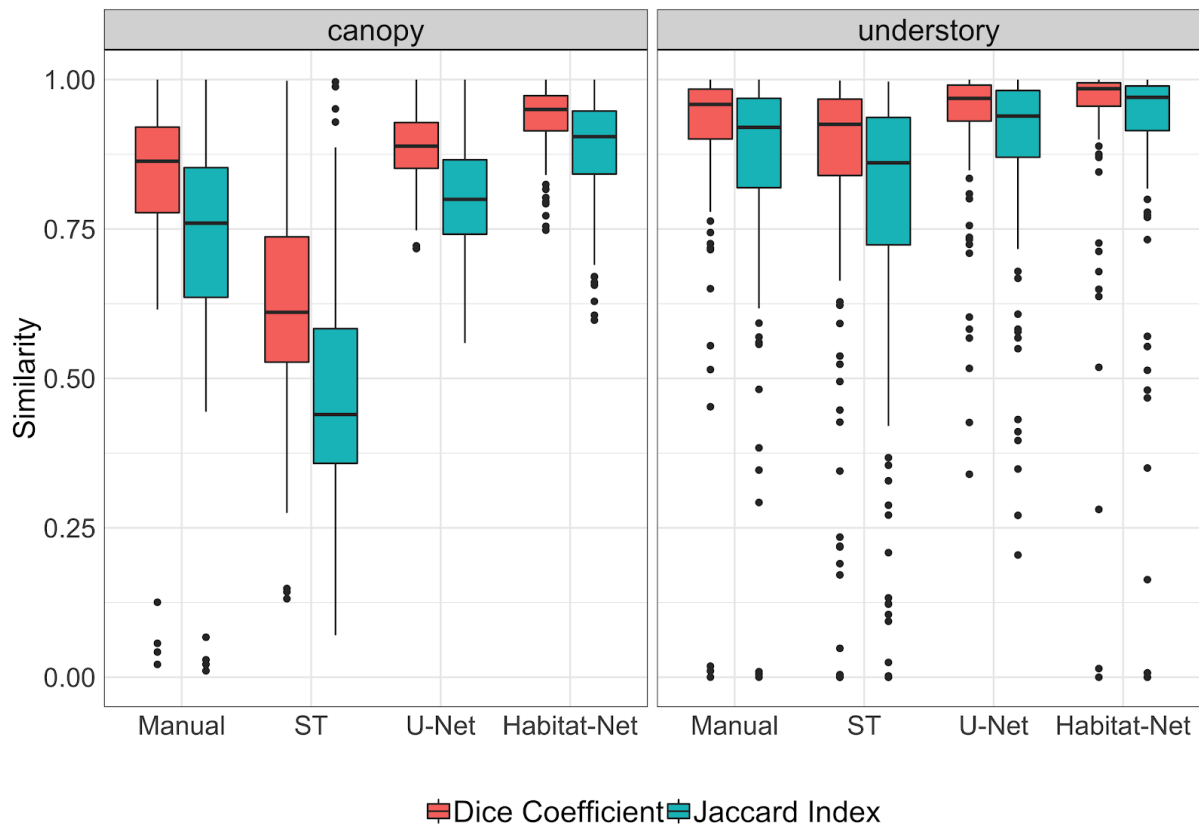
## A. Canopy Images



## B. Understory Images



553 **Figure 3:** Visual comparison of the quality of segmentation rasters predicted by Simple  
554 thresholding, U-Net, and Habitat-Net to the manually segmented reference images for  
555 different situations.



556 **Figure 4:** Box plots of the similarity scores (Dice coefficient and Jaccard index) between the  
557 image segmentation output by four methods and the reference manual segmentation for  
558 canopy and understory images.

559 **Tables**

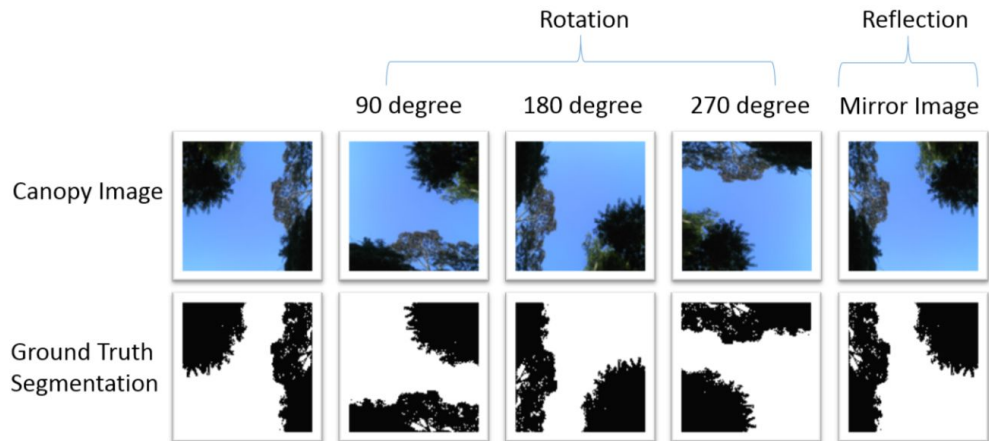
560 **Table 1:** Mean processing time per photo and similarity scores for the canopy dataset using  
 561 four methods.

| Method              | Processing time (seconds) | Dice coefficient |      |        | Jaccard index |      |        |
|---------------------|---------------------------|------------------|------|--------|---------------|------|--------|
|                     |                           | Mean             | SD   | Median | Mean          | SD   | Median |
| Manual              | 45                        | 0.84             | 0.16 | 0.86   | 0.74          | 0.19 | 0.76   |
| Simple thresholding | 0.103                     | 0.64             | 0.17 | 0.61   | 0.49          | 0.19 | 0.45   |
| U-Net               | 0.015                     | 0.89             | 0.07 | 0.89   | 0.81          | 0.11 | 0.80   |
| Habitat-Net         | 0.015                     | 0.94             | 0.06 | 0.95   | 0.88          | 0.09 | 0.90   |

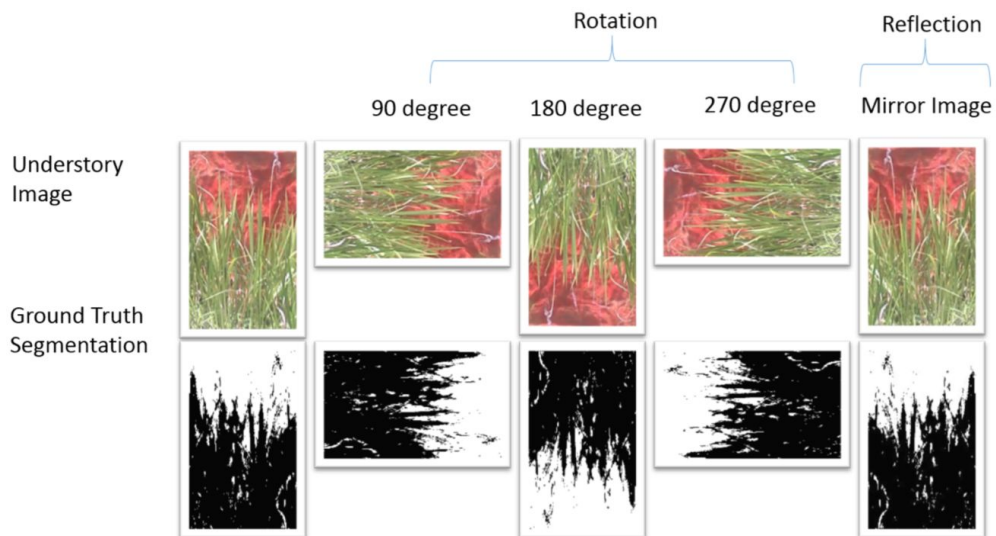
562 **Table 2:** Mean processing time per photo and similarity scores for the understory vegetation  
 563 dataset using four methods.

| Method              | Processing time (seconds) | Dice coefficient |      |        | Jaccard index |      |        |
|---------------------|---------------------------|------------------|------|--------|---------------|------|--------|
|                     |                           | Mean             | SD   | Median | Mean          | SD   | Median |
| Manual              | 65                        | 0.91             | 0.15 | 0.96   | 0.86          | 0.17 | 0.92   |
| Simple thresholding | 0.094                     | 0.83             | 0.25 | 0.93   | 0.77          | 0.27 | 0.86   |
| U-Net               | 0.015                     | 0.94             | 0.10 | 0.97   | 0.89          | 0.14 | 0.94   |
| Habitat-Net         | 0.015                     | 0.95             | 0.13 | 0.98   | 0.92          | 0.16 | 0.97   |

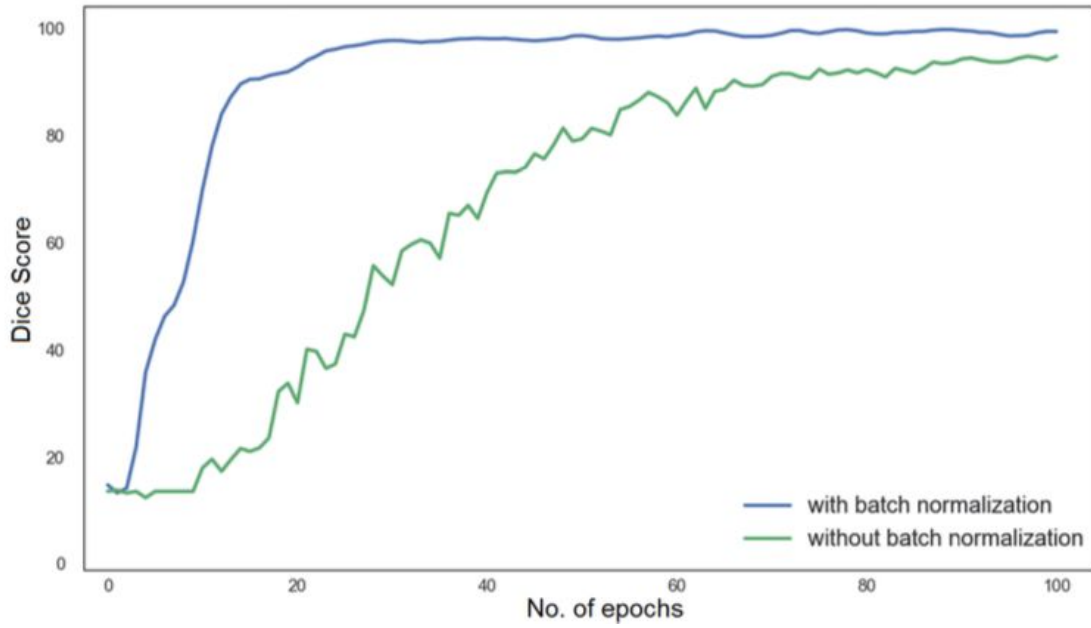
## Supplementary material



**Figure S1:** An example of data augmentation for the canopy image dataset



**Figure S2:** An example of data augmentation for the understory image dataset



**Figure S3:** Convergence with and without batch normalization.