# 1 DeepForest: A Python package for RGB deep

# 2 learning tree crown delineation

3 Ben. G. Weinstein[1], Sergio Marconi[1], Mélaine Aubry-Kientz[2], Gregoire Vincent[2], Henry

4 Senyondo[1], Ethan White[1]

5 [1]Department of Wildlife Ecology and Conservation, University of Florida, Gainesville,

6 Florida, USA

7 [2]AMAP, IRD, CNRS, INRA, Univ Montpellier, CIRAD, 34000 Montpellier, France

# 8 Abstract

9 1. Remote sensing of forested landscapes can transform the speed, scale, and cost

10 of forest research. The delineation of individual trees in remote sensing images is

11 an essential task in forest analysis. Here we introduce a new Python package,

12 DeepForest, that detects individual trees in high resolution RGB imagery using

13 deep learning.

14 2. While deep learning has proven highly effective in a range of computer vision

15 tasks, it requires large amounts of training data that are typically difficult to obtain

16 in ecological studies. DeepForest overcomes this limitation by including a model

17 pre-trained on over 30 million algorithmically generated crowns from 22 forests

18 and fine-tuned using 10,000 hand-labeled crowns from 6 forests.

19 3. The package supports the application of this general model to new data, fine

20 tuning the model to new datasets with user labeled crowns, training new models,

1

21     and evaluating model predictions. This simplifies the process of using and

22     retraining deep learning models for a range of forests, sensors, and spatial

23     resolutions.

24   4. We illustrate the workflow of DeepForest using data from the National Ecological

25     Observatory Network, a tropical forest in French Guiana, and street trees from

26     Portland, Oregon.

27   Keywords: Remote Sensing, Forests, Tree Crowns, Crown Delineation, NEON, Deep

28   learning, RGB

# Introduction

30   Airborne individual tree delineation is a central task for forest ecology and the

31   management of forested landscapes. The growth in sensor quality and data availability

32   has raised hopes that airborne tree maps can complement traditional ground-based

33   surveys (Hamraz et al. 2016). Most approaches to tree delineation in remote sensing

34   use three-dimensional LIDAR data (Coomes et al. 2017), which is currently available for

35   only a small fraction of the Earth's surface. In contrast, high resolution RGB data has

36   widespread coverage from commercial and government sources and is readily collected

37   using unmanned aerial vehicles. As a result, there is an increasing need for RGB-based

38   tree delineation approaches with easy to use open-source implementations.

39     The introduction of deep neural networks has greatly enhanced the performance

40   of remote sensing solutions for detecting objects in geospatial images (Zhu et al. 2017).

41   Deep learning models use a series of hierarchical layers to learn directly from training

42   data instead of using expert designed features. Initial layers learn general

2

43    representations, such as colors and shapes, and subsequent layers learn specific object

44    representations. There are several barriers to applying deep learning to ecological

45    applications including insufficient technical expertise, a lack of large amounts of training

46    data, and the need for significant computational resources. DeepForest provides easy

47    access to deep learning for tree delineation by creating a simple interface for training

48    object detection models, using them to make predictions, and evaluating the accuracy

49    of those predictions. DeepForest also includes a prebuilt model (based on Weinstein et

50    al. 2020) pre-trained on tens of millions of LiDAR generated crowns and fine-tuned

51    using over 10,000 hand-labeled crowns from diverse forests in the National Ecological

52    Observatory Network. Users can apply this model to detect trees in new imagery or

53    provide additional hand-labeled data to fine-tune performance for a specific site or forest

54    type. Predictions from the model for an average $1km^2$ tile can be made in 7 minutes on

55    a single CPU and DeepForest has built-in support for running on GPU resources to

56    dramatically increase the speed of prediction at large scales.

# 57  DeepForest Software

58    DeepForest is an open source (MIT license) Python package supporting Python 3.6 and

59    Python 3.7 and has been tested on Windows, macOS, and Linux operating systems. It

60    can be installed using the Python Package Index (https://pypi.org/project/deepforest/) or

61    using the conda package manager for Windows, Linux and OSX

62    (https://github.com/conda-forge/deepforest-feedstock). The software is openly

63    developed on GitHub (https://github.com/weecology/DeepForest) with automated

64    testing and each release is archived on Zenodo

3

65 (https://doi.org/10.5281/zenodo.2538143). All DeepForest functions are documented

66 online with reproducible examples (https://deepforest.readthedocs.io/) and video

67 tutorials.

## Prebuilt model

69 DeepForest currently includes one prebuilt model (available by running

70 `deepforest.use_release`) that was trained on data from the National Ecological

71 Observatory Network (NEON) using a semi-supervised approach outlined in Weinstein

72 et al. (2019, 2020) (Figure 1). The model was pretrained on data from 22 NEON sites

73 using an unsupervised LiDAR based algorithm (Silva et al. 2016) to generate millions of

74 moderate quality annotations for model pretraining. The pretrained model was then

75 retrained based on over 10,000 hand-annotations of RGB imagery from six sites

76 (MLBS, NIWO, OSBS, SJER, TEAK, LENO; see NEON site abbreviations S1). The full

77 workflow is shown in Figure 1. While LIDAR data is used to facilitate data generation for

78 the prebuilt model, prediction relies only on RGB data, allowing the model to be used to

79 detect trees using RGB imagery alone. This prebuilt model extends the methods from

80 Weinstein et al. (2019, 2020) by pretraining on a much larger number of trees (30

81 million LIDAR-generated crowns compared to 10 million in Weinstein et al. 2020) and

82 diversity of sites (22 instead of 4 in Weinstein et al. 2020). Additional details on the

83 modeling approach, data generation, and model evaluation are available in Weinstein et

84 al (2019, 2020) and a brief summary is provided in S2. This model can be used directly

85 to make predictions for new data or used as a foundation for retraining the model using

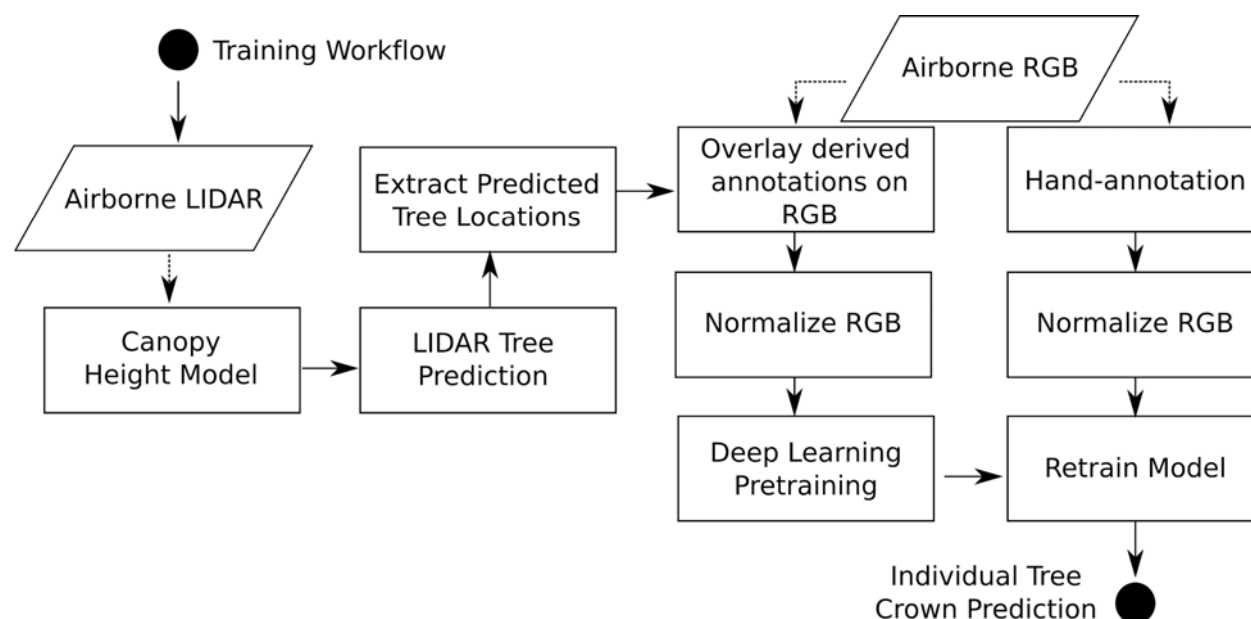86 labeled data from a new application.

4

Figure 1. Prebuilt model training workflow. Redrawn from Weinstein et al. (2020). Parallelograms in the workflow indicate input data, rectangles indicate an algorithmic step, and circles indicate the start and end of the workflow. The two sub-flows on the right-side of the figure can run in parallel and outline the pre-training and fine-tuning stages of the overall model fitting process.

## Training

Tree crown delineation is a challenging task and a single model cannot realistically be expected to capture the tremendous taxonomic diversity at a global scale. This means that to perform optimal crown delineation for a particular forest requires training or fine-tuning using data from a local area. A key advantage of DeepForest's neural network structure is that users can retrain the prebuilt model to learn new tree features and image backgrounds while leveraging information from the existing model weights based on data from a diverse set of forests. Fine-tuning neural networks starting from an initial model requires less training data to produce reasonable results (Shin et al. 2016).

5

102  Known as "transfer learning", this ability is important because training deep learning

103  models from scratch often requires tens of thousands of labeled data points for

104  ecological tasks (Weinstein 2018). In contrast, fine-tuning the prebuilt model with as few

105  as 1000 hand labeled trees can provide significant improvement and be accomplished

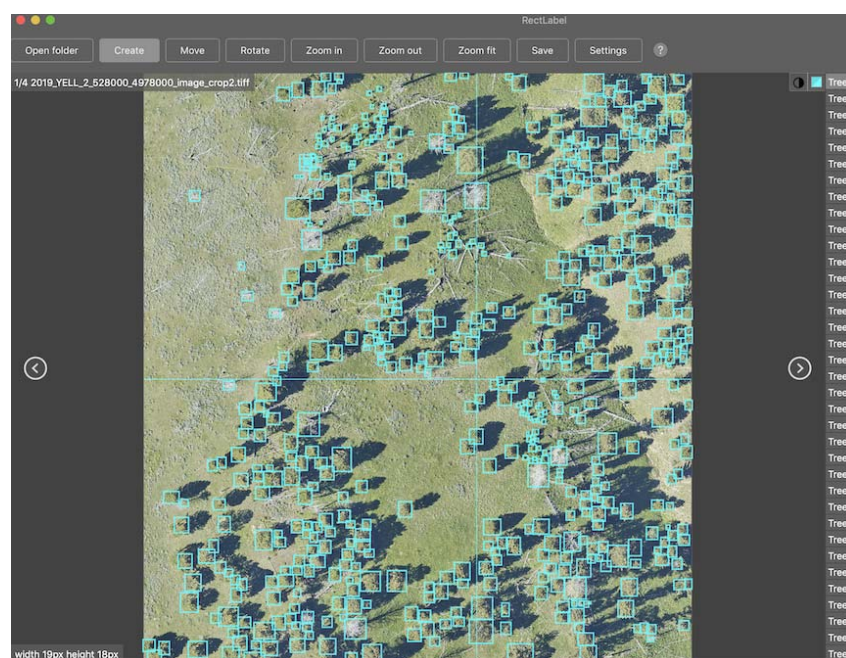106  in approximately 8-10 hours (Weinstein et al. 2020).

107



108  Figure 2. Screenshot of hand-annotated RGB image from NEON site YELL near Frog

109  Rock, WY. For optimal training, all crowns in an image should be annotated.

110      The standard training process starts with generating local training data by hand-

111  labeling trees in images by placing a bounding box around each visible tree (Figure 2).

112  This can be done using either image labeling tools (e.g., RectLabel) or GIS software

113  (e.g., ArcGIS, QGIS) and DeepForest includes helper functions to convert common

114  formats (XML and shapefiles) into a csv format. Annotations can be made on images of

115  any size, but training the model requires images with fixed standard dimensions. The

116  prebuilt model was trained on square crops of length 400px (40m at 0.1m resolution),

6

117   which provides a good balance between image size and providing the model landscape

118   context for prediction. DeepForest includes a `preprocess.split_raster` function

119   that creates a set of appropriately sized images for training using a sliding window

120   approach. The size of these input windows are optimized for the 10cm data used in

121   training the prebuilt model. The upper resolution limit for tree crown delineation is

122   currently unknown, as well as the optimal size of the input windows when performing

123   predictions at coarser scales.

124       Training can be performed by fine-tuning the prebuilt model or training only using

125   the local training data (using `deepforest.train`). Training deep learning models

126   requires a number of parameter choices such as batch size and number of epochs. For

127   users less familiar with training deep learning models, DeepForest comes with a

128   standard configuration file with reasonable defaults. While some parameter exploration

129   will always be helpful, our aim is to make these innovations available even to novice

130   users. Optional GPU support and model customization allow more experienced users to

131   quickly develop and test larger and more complex models. Data augmentation to

132   randomly crop and flip training images is also supported. This strategy is often useful to

133   reduce overfitting when training on small datasets (Zoph et al. 2019) but has not been

134   extensively tested for tree crown delineation. For additional recommendations for

135   optimal model training see the online documentation (https://deepforest.readthedocs.io/)

136   and Appendix S2.

## Evaluation

138   Deep neural networks have millions of parameters and can readily overfit, producing

139   high scores on training data, while performing poorly on new images. This makes it

7

140   essential to evaluate performance on held-out test data. To evaluate a set of

141   annotations, users follow the same pattern as with training data: 1) Annotate one or

142   more images of trees; 2) Cut the images into smaller windows for evaluation; and 3)

143   Format annotations into a csv file using DeepForest's utility functions. The

144   `deepforest.evaluate_generator` method can then be used to evaluate the

145   performance of the predictions for this test data using the mean average precision

146   (mAP). mAP combines precision and recall into a single metric measuring the area

147   under the precision-recall curve resulting in a score ranging from 0 to 1. In our

148   experience, mAP scores above 0.5 are usable for scientific application, but the

149   appropriate value depends on the particular research goal and application.

150   ## Prediction

151   After a model has been trained and evaluated, it can be applied to a larger collection of

152   images to estimate the locations of trees at larger scales. High resolution images

153   covering wide geographic extents cannot fit into memory during prediction and would

154   yield poor results due to the size and density of bounding boxes. DeepForest has a

155   `deepforest.predict_tile` method for automating the process of splitting the tile

156   into smaller overlapping windows, performing prediction on each of the windows, and

157   then reassembling the resulting annotations. Each bounding box annotation is returned

158   with its xmin, ymin, xmax, ymax coordinates, and predicted probability score (the

159   probability that the bounding box represents a tree) ranging from 0-1, with higher values

160   indicating greater confidence in the prediction. To reduce overcounting among

161   overlapping tiles, DeepForest sorts predictions by confidence scores and removes lower

162   scoring overlapping boxes (i.e., non-max suppression).
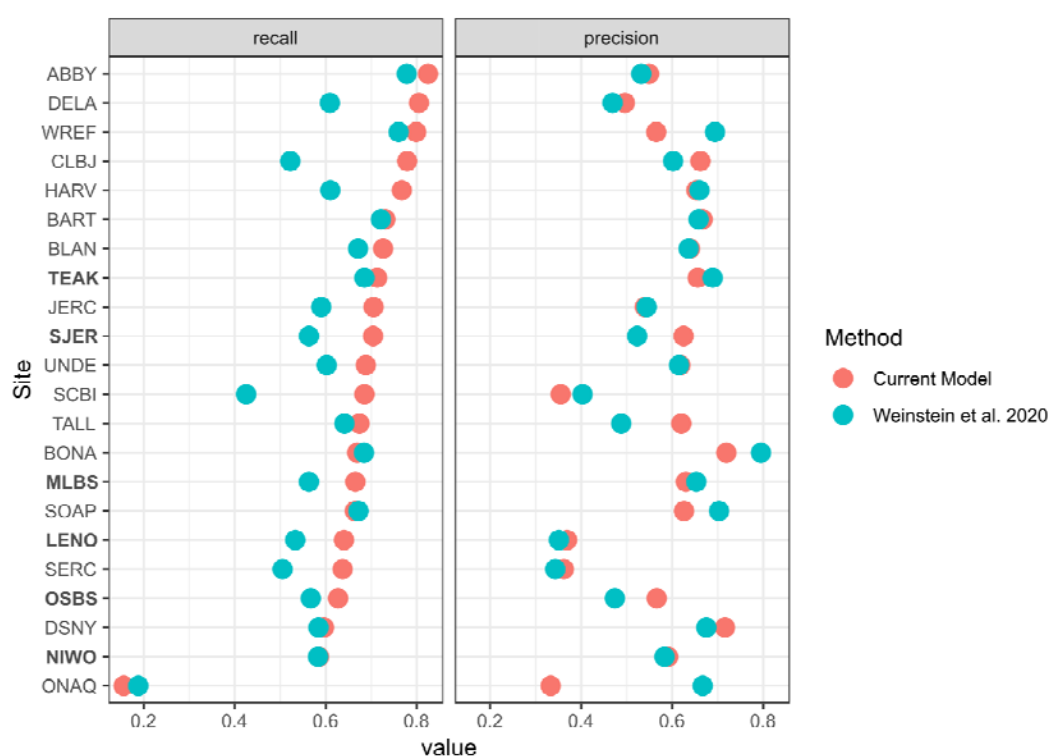
8

# Case Studies

## National Ecological Observatory Network

To evaluate model performance across a range of forest types, we used data from the

National Ecological Observatory Network to predict crowns in sites across the United

States. This dataset consists of 212 images containing 5852 trees from 22 sites that is

part of an upcoming tree crown benchmark data package (Weinstein et al. 2020).

Training and evaluation data are separated by at least 1 km when they occur at the

same site. Evaluation data were created by viewing RGB images and manually

delineating tree crown boxes for all visible trees. Annotations were cross-referenced

with field collected positions of tree stems (from the NEON Vegetation Structure

dataset; NEON ID: DP1.10098.001) within each plot when available. Following

Weinstein et al. (2019, 2020), we used precision, defined as the fraction of predicted

crowns match real trees, and recall, defined as the fraction of all evaluation trees that

are correctly detected for evaluation. Following the standard evaluation for object

detection in the computer vision literature (Ren et al. 2015), we considered predictions

with Intersection over Union (IoU) scores of 0.5 as true positives. IoU, also known as

the Jaccard Index, is the area of intersection between the prediction and evaluation

crown, divided by the joint area of the combined prediction and evaluation crowns. We

assessed the performance of the prebuilt model at all 22 NEON sites and also

compared the performance to a previous version of this model (Weinstein et al. 2020)

that was only trained on data from 4 NEON sites.

9

184      Across all sites the average recall per image for the prebuilt model was 72% and

185    the precision was 64%. Model performance varies across NEON sites, but most sites

186    have both precision and recall values greater than 50% (Figure 3). The model performs

187    similarly regardless of whether there was hand-annotated training data from the same

188    site (Figure 3). Visual assessment of predictions across forest types reveals good

189    overall correspondence between predicted bounding boxes and observations, with most

190    errors resulting from insufficient overlap between observed and predicted tree crowns,

191    rather than the model missing a tree entirely (Figure 4). The prebuilt model used by

192    DeepForest was fit to data from 22 NEON sites and outperforms the previous 4 site

193    model (Weinstein et al. 2020) at 19 of 22 sites for recall and 16 of 22 sites for precision,

194    demonstrating that increasing the diversity and amount of training data has improved

195    the performance of the model. These results demonstrate that the prebuilt model can

196    make reasonable predictions in forests ranging from deciduous forests of the Northeast,

197    to southern pinelands, to coniferous forests of the mountain west.

198      The site with the worst performance is Onaqui, Utah (ONAQ), which is a desert

199    scrub site with a different vegetation structure from any of the training data. The site is

200    almost treeless and includes trees with short and gnarled stature. This highlights the

201    importance of using local training data to reduce uncertainty when working with data

202    that is not well represented in the training data for the prebuilt model. In these contexts,

203    the value of the prebuilt model is that it reduces the needed training sizes when applied

204    to new conditions. This has the potential to support training with small amounts of data

205    for applications to a wide array of questions surrounding tree health and ecology. For

206    example, training a model specific to bare trees could allow studies of broad-scale pest

10

207 outbreaks or timing of deciduous phenology. Initial tests at the Soaproot Saddle, CA site

208 ('SOAP' in Figure 4, 3$^{rd}$ row) show the prebuilt model can detect standing dead trees

209 when visible. Adding additional training data could allow broad scale analysis of tree

210 health when comparing images across time.



211

212 Figure 3. Precision and recall scores for hand-labeled evaluation images from the

213 National Ecological Observatory Network (current prebuilt model in red, Weinstein et al.

214 2020 in blue). Sites in bold had hand-labeled data included in training the current

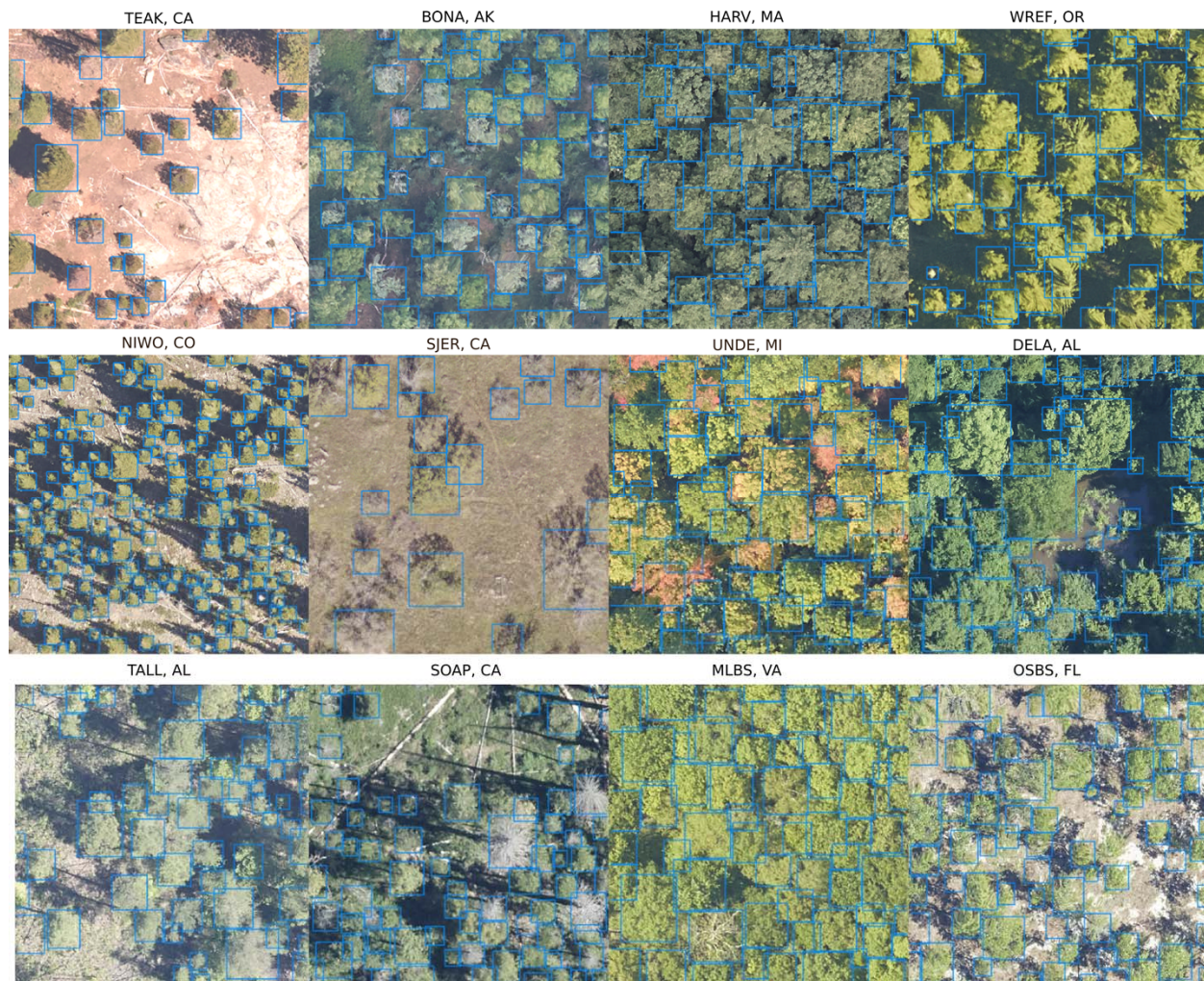215 prebuilt model. See S1 for site abbreviations.

11

216

217 Figure 4. Panel of tree predictions from a broad range of evaluation images in the

218 National Ecological Observatory Network with predicted tree crown boxes are in blue.

219 Each image is labeled with the NEON site abbreviation and state. See S1 for site

220 abbreviations.

## French Guiana Tropical Forest

222 The DeepForest prebuilt model was trained on data from the United States that was

223 collected using fixed-winged aircraft at 10cm resolution and provided as $1km^2$

224 orthomosaics. Therefore, two key questions are: 1) Does this model generalize to

12

225 images collected in new locations or using different acquisition hardware; and 2) how

226 useful are the (re)training features of the software for improving performance in novel

227 contexts? It is also important to understand how the DeepForest RGB model compares

228 to LiDAR-based models from recently published work.

229     To address these questions, we used data from a recently published competition

230 comparing LIDAR tree segmentation algorithms using remote sensing from French

231 Guiana (Aubry-Kentz et al. 2019). In the original competition, each team was sent

232 unlabeled data to predict and the evaluation data was kept private. This process was

233 repeated for this paper, with the third author (Aubry-Kientz) running evaluation scores

234 for the DeepForest predictions made by the corresponding author (Weinstein).

235 Predictions were run on a Mac laptop with a 3.1 GHz Intel Core i5 processor.

236 Predictions from each algorithm were compared to hand-delineated evaluation crowns

237 based on field observation and manual comparison with RGB and LiDAR data.

238 Validation crowns were delineated as polygons, rather than the rectangular bounding

239 boxes generated by DeepForest. This case study also provides information on whether

240 DeepForest's approach of predicting rectangular bounding boxes leads to lower

241 prediction accuracy than methods producing polygons. Algorithm recall was scored

242 based on the proportion of labeled trees predicted with IoU scores of greater than 0.5.

243 Precision was not calculated because not all crowns in the test imagery were delineated

244 (see Figure 5).

13

245



246

Figure 5. RGB images collected over a tropical forest in French Guiana and example of

manually segmented crowns used to evaluate the segmentation.

We used DeepForest to detect tree crowns using three approaches 1) the

prebuilt DeepForest model with no local training; 2) a model fit solely to 5018 local

hand-annotated crowns (annotated by BW using only the RGB data on tiles separate

from the evaluation data)  and 3) the prebuilt model fine-tuned using the local

annotations. RGB tiles were divided into 800px windows for model training and

evaluation. The default patch size of 400px was increased to 800px to minimize the

edge effect of overlapping crowns. Models 2 and 3 were trained for 7 epochs with a

runtime of approximately 11 minutes/CPU on a laptop, which demonstrates that while

advanced GPU hardware is convenient for training large datasets, fine-tuning and

training on small datasets can be done locally on CPU.

The prebuilt model performed well on this novel data with a recall of 0.64, close

to the 0.71 recall for the best performing LIDAR based algorithm from Aubry-Kientz et

al. (2019). Training only on the 5018 local annotations resulted in a poorer recall of

0.35. Retraining the prebuilt model with the local annotations produced the best results

14

263   with a recall of 0.78, slightly better than the highest performing LiDAR algorithm from

264   Aubry-Kientz et al. (2019). This analysis is not sufficient to draw general conclusions

265   about RGB versus LiDAR-based methods, but these results do suggest that

266   DeepForest is competitive with state-of-the-art LiDAR-based approaches. Overall, the

267   case study demonstrates the utility of DeepForest both using the prebuilt model and

268   using local retraining to improve crown delineation based on local conditions.
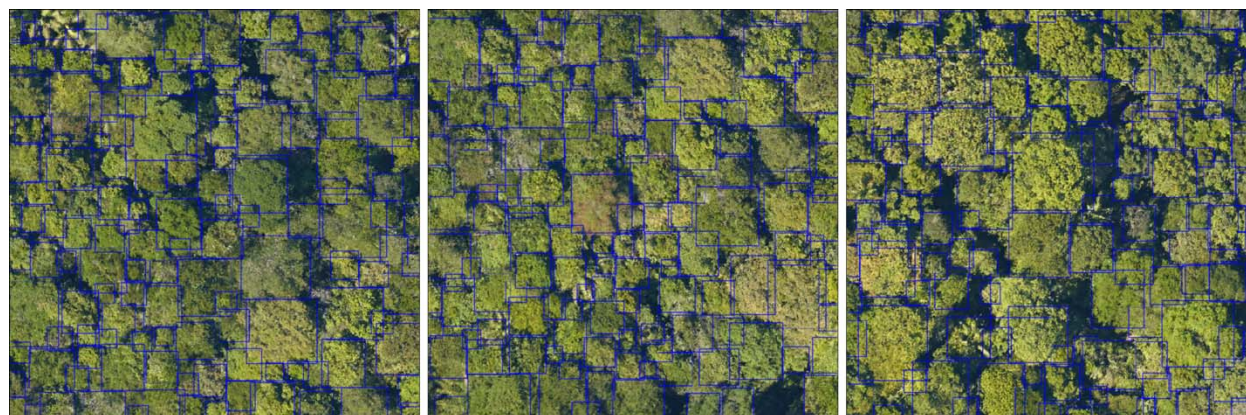
269



270   Figure 6. Predictions made on a tropical forest in French Guiana using the prebuilt

271   model retrained with local annotations. Each individual tree is labeled with a blue

272   bounding box.

## Portland Street Trees

274   DeepForest's use of widely available RGB data provides the potential for it to be used

275   across very large spatial extents. Scaling up is challenging because algorithms need to

276   handle large ranges of habitat types and because the resolution of the data available

277   over large areas is typically coarser. To explore how DeepForest performs using

278   coarser resolution data in unique habitats, we applied both the prebuilt model and a

279   retrained model to crown delineation of street trees in an urban environment. The

280   locations of urban trees are important for ecological, sociological and public

15

281    infrastructure applications. In addition, the urban environment is very different from the

282    natural environments on which the prebuilt model was trained. The image data from the

283    Oregon Statewide Imagery Program is also coarser at 0.3m spatial resolution (1ft), a

284    resolution that is widely available as part of the National Agriculture Imagery Program

285    (NAIP - https://www.fsa.usda.gov/programs-and-services/aerial-photography/imagery-

286    programs/naip-imagery/).

287        We used imagery from the Portland metro area that overlapped with the Portland

288    Street Trees dataset (http://gis-pdx.opendata.arcgis.com/datasets/street-trees). The

289    street trees dataset contains geospatial information for the majority of trees accessible

290    from public roads in the metro area. Not all trees in an image are labeled, since many

291    trees occur on private property and are not mapped. We divided the RGB imagery into

292    geographically distinct training and test datasets and used the street trees dataset to

293    guide hand-annotation of a small number of tree crowns (n=1033). Annotation by hand

294    took approximately three hours and covered a small geographic area of mixed urban

295    development, empty lots and ballfields (Figure 6). The street trees data was collected

296    prior to the RGB images and was cleaned to remove trees that had been cut down or

297    were obvious errors (e.g. trees located in the middle of buildings). To evaluate the street

298    tree case study, we used field collected location of the tree stems to measure tree recall

299    and the rate of undersegmentation. Recall was defined as the proportion of street tree

300    locations that were contained within a predicted tree bounding box. Undersegmentation

301    rate was defined as the proportion of predicted boxes that matched more than one

302    street tree. Minimizing undersegmentation is challenging because trees growing close

16

303    together can appear to be a single tree from above and is therefore best evaluated

304    against ground collected data.

305         We found that evaluating and retraining on data with coarser resolution than the

306    prebuilt model required careful choosing of the size of the focal view. The prebuilt model

307    was originally trained on a 40m focal view (400px windows with 0.1m data). Data

308    exploration on the coarser data source showed that larger focal views of 60-120m

309    performed better than maintaining the original 40m view, and 60m was chosen for this

310    analysis. In general, we expect that the focal view size should increase with coarser

311    resolution data, but this remains an area of further exploration.

312         As with the tropical forest case study, we found that the prebuilt model performed

313    reasonably well (recall = 0.55; undersegmentation = 0.25) and retraining with a small

314    amount of local training data significantly improved algorithm performance with an

315    increase in recall and decrease in undersegmentation (recall = 0.72; undersegmentation

316    = 0.17; Figure 7). Visual inspection shows that many of the errors in using the retrained

317    model are for small trees difficult to resolve in the imagery, or tree types not present in

318    the limited training data (e.g. ornamental trees with a deep purple hue).

17

Figure 7. Predictions for the Portland street tree case study. Bounding box predictions from the prebuilt model are in orange. Bounding box predictions from the retrained model using local data are in blue. Street tree locations are marked in purple.

# Conclusion

DeepForest provides an open source software package for: 1) delineating tree crowns in RGB imagery, 2) evaluating the performance of that crown delineation using hand

18

327    labeled evaluation data, and 3) training new models and fine-tuning of the included

328    prebuilt model to support prediction tailored to specific forest types. The inclusion of a

329    prebuilt model allows users to benefit from the strengths of deep learning without

330    needing to deal with many of the challenges. Given the enormous diversity of tree

331    appearance at a global scale, defining a single unified model for tree crown delineation

332    is challenging. To address this, DeepForest provides an explicit retraining method to

333    improve performance for specific use cases. This allows the user to decide what level of

334    accuracy is required for the target question, and then annotate local data and retrain the

335    model to produce predictions with sufficient accuracy for their use case. We recommend

336    defining a clear evaluation dataset, setting a threshold for desired performance before

337    training, and using evaluation data that is geographically separate from the training data

338    to ensure that the prediction threshold holds outside of the training region.

339        The minimal spatial resolution for accurate tree prediction using this software

340    remains unknown and may ultimately relate to the desired ecological or management

341    question. Analysis of the NEON data show that individual tree segmentation is

342    achievable at 10cm.  The Portland Street trees example shows that 30 cm data (which

343    is publicly available for many states and counties) provides reasonable delineations.

344    However, the accuracy will not be as high as with higher resolution data, and further

345    analysis at this resolution is necessary. One meter resolution imagery is increasingly

346    available at near continental scales (e.g., NAIP 1m imagery which provides nearly

347    complete coverage of the United States; https://www.fsa.usda.gov/programs-and-

348    services/aerial-photography/imagery-programs/naip-imagery/). It is unlikely that these

19

349    data will be effective at distinguishing small individual trees, but it may be useful in

350    identifying large trees or clusters of trees in sparse landscapes.

351         To support the broad application of predictions from DeepForest, these

352    predictions can be easily exported for use in further analysis and combination with other

353    sensor products for forest research. Individual tree crown delineation is often the first

354    step in key remote sensing analyses of forested landscapes, including biomass

355    estimation (Kamoske et al. 2019), species classification (Maschler et al. 2018), and leaf-

356    trait analysis (Marconi et al. 2019). DeepForest both ingests and outputs crowns in an

357    easily accessible, standardized annotation format, and will facilitate further

358    improvements in the prebuilt model based on community contributions.

359    # Data Availability

360    DeepForest source code is available on GitHub

361    (https://github.com/weecology/DeepForest) and archived on Zenodo

362    (https://doi.org/10.5281/zenodo.2538143). The code for the case studies is available in

363    a separate repo (https://github.com/weecology/DeepForest_demos).  The in-

364    development  version of the NEONTreeEvaluation benchmark is available online

365    (https://github.com/weecology/NeonTreeEvaluation) and will continue to be updated as

366    more images are annotated. The Oregon RGB imagery was provided by the Oregon

367    Statewide Imagery Program 2018: 1 foot orthophotography of western Oregon: State of

368    Oregon data release, https://www.oregon.gov/geo/Pages/imagery_data.aspx

20

# Author Contributions

370 BW, SM, EW conceived of the project, designed the package and wrote the manuscript.

371 MAK, GV, and HS performed analysis, provided package improvements and edited the

372 manuscript.

# Literature Cited

382 Aubry-Kientz, M., Dutrieux, R., Ferraz, A., Saatchi, S., Hamraz, H., Williams, J.,

383 Coomes, D., Piboule, A., Vincent, G., 2019. A Comparative Assessment of the

384 Performance of Individual Tree Crowns Delineation Algorithms from ALS Data in

385 Tropical Forests. Remote Sens. 11, 1086. https://doi.org/10.3390/rs11091086

386 Coomes, D.A., Dalponte, M., Jucker, T., Asner, G.P., Banin, L.F., Burslem, D.F.R.P.,

387 Lewis, S.L., Nilus, R., Phillips, O.L., Phua, M.H., Qie, L., 2017. Area-based vs tree-

21

388   centric approaches to mapping forest carbon in Southeast Asian forests from airborne

389   laser scanning data. Remote Sens. Environ. 194, 77–88.

390   https://doi.org/10.1016/j.rse.2017.03.017

391   Dalponte, M., Coomes, D.A., 2016. Tree-centric mapping of forest carbon density from

392   airborne laser scanning and hyperspectral data. Methods Ecol. Evol. 7, 1236–1245.

393   https://doi.org/10.1111/2041-210X.12575

394   Dalponte, M., Frizzera, L., Gianelle, D., 2018. How to map forest structure from aircraft ,

395   one tree at a time. Ecol. Evol. 8, 5611–5618. https://doi.org/10.1002/ece3.4089

396   Ferraz, A., Saatchi, S., Mallet, C., Meyer, V., 2016. Lidar detection of individual tree size

397   in tropical forests. Remote Sens. Environ. 183, 318–333.

398   https://doi.org/10.1016/j.rse.2016.05.028

399   Hamraz, H., Contreras, M.A., Zhang, J., 2016. A robust approach for tree segmentation

400   in deciduous forests using small-footprint airborne LiDAR data. Int. J. Appl. Earth Obs.

401   Geoinf. 52, 532–541. https://doi.org/10.1016/j.jag.2016.07.006

402   Kamoske, A.G., Dahlin, K.M., Stark, S.C., Serbin, S.P., 2019. Leaf area density from

403   airborne LiDAR: Comparing sensors and resolutions in a temperate broadleaf forest

404   ecosystem. For. Ecol. Manage. 433, 364–375.

405   https://doi.org/10.1016/j.foreco.2018.11.017

406   Marconi, S. Graves, S. Weinstein, B. Bohlman, S. White, E. 2020. Rethinking the

407   fundamental unit of ecological remote sensing: Estimating individual level plant traits at

408   scale. bioRxiv 556472; doi: https://doi.org/10.1101/556472

22

409  Maschler, J., Atzberger, C., Immitzer, M., 2018. Individual Tree Crown Segmentation

410  and Classification of 13 Tree Species Using Airborne Hyperspectral Data. Remote

411  Sens. 10, 1218. https://doi.org/10.3390/rs10081218

412  Ren, S., He, K., Girshick, R., Sun, J., 2015. Faster r-cnn: Towards real-time object

413  detection with region proposal networks. Nips 91–99.

414  https://doi.org/10.1109/TPAMI.2016.2577031

415  Shin, H., Roth, H.R., Gao, M., Lu, L., Member, S., Xu, Z., Nogues, I., Yao, J., Mollura,

416  D., Summers, R.M., 2016. Deep Convolutional Neural Networks for Computer-Aided

417  Detection□: CNN Architectures, Dataset Characteristics and Transfer Learning 35,

418  1285–1298.

419  Weinstein, B.G., 2018. A computer vision for animal ecology. J. Anim. Ecol. 87, 533–

420  545. https://doi.org/10.1111/1365-2656.12780

421  Weinstein, B.G., Marconi, S., Bohlman, S., Zare, A., White, E., 2019. Individual Tree-

422  Crown Detection in RGB Imagery Using Semi-Supervised Deep Learning Neural

423  Networks. Remote Sens. 11, 1309. https://doi.org/10.3390/rs11111309

424  Weinstein, B.G., Marconi, S., Bohlman, S.A., Zare, A., White, E.P., 2020. Cross-site

425  learning in deep learning RGB tree crown detection. Ecol. Inform. 56, 101061.

426  https://doi.org/10.1016/j.ecoinf.2020.101061

427  Weinstein, B. (2020, March 21). weecology/NeonTreeEvaluation: Zenodo Release

428  (Version 1.1). Zenodo. http://doi.org/10.5281/zenodo.3723357

23

429    Zhu, X.X., Tuia, D., Mou, L., Xia, G.-S., Zhang, L., Xu, F., Fraundorfer, F., 2017. Deep

430    Learning in Remote Sensing: A Comprehensive Review and List of Resources. IEEE

431    Geosci. Remote Sens. Mag. 5, 8–36. https://doi.org/10.1109/MGRS.2017.2762307


432    Zoph, B., Cubuk, E. D., Ghiasi, G., Lin, T. Y., Shlens, J., & Le, Q. V., 2019. Learning

433    data augmentation strategies for object detection. arXiv preprint arXiv:1906.11172.

434    Supplementary 1

435    NEON Site Abbreviations

| Site Name | Site ID | State | Latitude | Longitude |
|---|---|---|---|---|
| *Onaqui* | ONAQ | UT | 40.17759 | -112.45244 |
| *Santa Rita Experimental Range* | SRER | AZ | 31.91068 | -110.83549 |
| *Niwot Ridge Mountain Research Station* | NIWO | CO | 40.05425 | -105.58237 |
| *Yellowstone Northern Range (Frog Rock)* | YELL | WY | 44.95348 | -110.53914 |
| *Wind River Experimental Forest* | WREF | WA | 45.82049 | -121.95191 |
| *Abby Road* | ABBY | WA | 45.76243 | -122.33033 |
| *Lower Teakettle* | TEAK | CA | 37.00583 | -119.00602 |
| *Caribou-Poker Creeks Research Watershed* | BONA | AK | 65.15401 | -147.50258 |
| *Soaproot Saddle* | SOAP | CA | 37.03337 | -119.26219 |
| *San Joaquin Experimental Range* | SJER | CA | 37.10878 | -119.73228 |
| *Lenoir Landing* | LENO | AL | 31.85388 | -88.16122 |
| *UNDERC* | UNDE | MI | 46.23388 | -89.53725 |
| *Ordway-Swisher Biological Station* | OSBS | FL | 29.68927 | -81.99343 |
| *Disney Wilderness Preserve* | DSNY | FL | 28.12504 | -81.4362 |
| *Jones Ecological Research Center* | JERC | GA | 31.19484 | -84.46861 |
| *Blandy Experimental Farm* | BLAN | VA | 39.06026 | -78.07164 |
| *Smithsonian Environmental Research Center* | SERC | MD | 38.89008 | -76.56001 |
| *Smithsonian Conservation Biology Institute* | SCBI | VA | 38.89292 | -78.1395 |
| *Bartlett Experimental Forest* | BART | NH | 44.06388 | -71.28731 |
| *Konza Prairie Biological Station* | KONZ | KS | 39.10077 | -96.56309 |
| *Talladega National Forest* | TALL | AL | 32.95046 | -87.39327 |
| *Dead Lake* | DELA | AL | 32.54172 | -87.80389 |
| *Mountain Lake Biological Station* | MLBS | VA | 37.37828 | -80.52484 |
| *The University of Kansas Field Station* | UKFS | KS | 39.04043 | -95.19215 |
| *Harvard Forest* | HARV | MA | 42.5369 | -72.17266 |

436

25

437  S2: Technical background to DeepForest and advice on model development

438  The goal of this paper is to provide a robust open-source implementation for RGB tree crown

439  delineation. For detailed information on data generation and model testing see Weinstein et al.

440  (2019, 2020). Here we provide an overview with recommendations for parameters that could

441  affect model performance. For more discussion with code examples see

442  https://deepforest.readthedocs.io/ .

443  DeepForest is a deep learning object detection model. Deep learning uses a series of

444  hierarchical filters to connect low-level image features, such as colors and shape, to high-level

445  concepts such as crown contours. Deepforest uses the keras-retinanet implementation (Gaiser

446  et al 2017) of the Retinanet one-stage object detector with a Resnet-50 classification backbone

447  (Ren et al. 2015). This classification backbone was pretrained on the ImageNet dataset (He et

448  al., 2016). This pretraining is useful for reducing training time and minimize overfitting to smaller

449  ecological datasets. Given an input image, the model predicts both the location of bounding

450  boxes and object classes simultaneously. In the DeepForest implementation, the prebuilt model

451  has only 1 "Tree" class. All instances of the target class are predicted, such that if there are 40

452  trees in the image, an ideal model will predict 40 bounding boxes. Deep learning models have

453  millions of parameters and take significant computational resources. To reduce memory

454  consumption, we cut each prediction image into 40m by 40m windows with an overlap of 5%.

455  Using a pool of unsupervised LiDAR-based tree predictions generated using Silva et al (2016),

456  we pretrained the network with a batch size of 20 on 2 Tesla K80 GPU for 5 epochs. To align

457  these unsupervised classifications with the ImageNet pretraining weights, we normalized the

458  RGB channels by subtracting the ImageNet mean from each channel. We then retrained the

459  network using the hand-annotated data for 40 epochs. Data augmentation of random flips and

460  translations was tested and found to have little effect on the final score.

461  Key parameter choices

26

462 Given the diversity of forests, data resolutions and image backgrounds, it is unlikely that a

463 universal tree crown method exists without some parameter tuning. The prebuilt model does a

464 good job in many situations, but to get the most of of DeepForest, some parameter exploration

465 will be needed. Using the prebuilt model, users should consider:

466 • Patch Size

467 The prebuilt model was built on 40m crops with 0.1m imagery (400^2 pixels). The model can

468 generalize to new resolutions, (see street trees example), but may need to vary the image patch

469 size. For example, on the 30cm NAIP data, we found that increasing patch size from 40 to 60

470 meters per image improved model performance. While an exact equation between image

471 resolution and optimal patch size is unknown, we anticipate that coarser resolution data require

472 larger patches to gather more context around visible tree crowns.

473 • IoU suppression threshold

474 Object detection models have no inherent logic about overlapping bounding boxes. For tree

475 crown detection, we expect trees in dense forests to have some overlap, but not be completely

476 intertwined. We therefore apply a postprocessing filter called 'non-max-suppression', which is a

477 common approach in the broader computer vision literature. This routine starts with the boxes

478 with the highest confidence scores and removes any overlapping boxes greater than the

479 intersection-over-union threshold (IoU). Intersection-over-union is the most common object

480 detection metric, defined as the area of intersection between two boxes divided by the area of

481 union. If users find that there is too much overlap among boxes, increasing the IoU will return

482 fewer boxes. To increase the number of overlapping boxes, reduce the IoU threshold.

483

484 Training New Models

27

485     Ultimately, the best performance will come from training local data. Even a small amount of local

486     training data can be beneficial. When training images, it is important to label all visible trees, not

487     just the trees that the user has high confidence in. Skipping trees in an image will teach the

488     model to skip trees during prediction. During annotation, we highly suggest users spend the

489     time to answer 2 questions: 1) What kind of data am I trying to predict? Capturing the variability

490     and the broad range of tree taxonomy and presentation will make development go more

491     smoothly. 2) What kind of accuracy do I need to answer my question? It is natural to want the

492     best model possible, but one can waste a time trying to eek out another 5% of recall without

493     understanding whether that increase in performance will improve our understanding of a given

494     ecological or natural resource question.

495        • Batch size

496     Neural networks are often trained in batches of images, since the entire dataset is often too

497     large to read into memory at once. The size of these batches affects both the speed of training

498     (larger batches train faster) and the stability of training (larger batches lead to more consistent

499     results). The default batch size of 1 is chosen because it is not possible to anticipate the

500     available memory and should be increased whenever possible. Typically, batch sizes are evenly

501     divisible by the size of the entire dataset.

502        • Epochs

503     An 'epoch' is one iteration of training on all images. The more times the model sees the training

504     data, the higher the training accuracy. However, deep learning models have millions of

505     parameters and can easily overfit. For small datasets (< 10,000 annotations), 5-10 epochs

506     should be sufficient. If training for longer, make sure to carefully evaluate validation accuracy on

507     out-of-sample data.

508

509

28