

**A peer-reviewed version of this preprint was published in PeerJ on 28 February 2019.**

[View the peer-reviewed version](https://peerj.com/articles/5843) (peerj.com/articles/5843), which is the preferred citable publication unless you specifically need to cite this preprint.

Marconi S, Graves SJ, Gong D, Nia MS, Le Bras M, Dorr BJ, Fontana P, Gearhart J, Greenberg C, Harris DJ, Kumar SA, Nishant A, Prarabdh J, Rege SU, Bohlman SA, White EP, Wang DZ. 2019. A data science challenge for converting airborne remote sensing data into ecological information. PeerJ 6:e5843 <https://doi.org/10.7717/peerj.5843>

# A data science challenge for converting airborne remote sensing data into ecological information

**Sergio Marconi**<sup>1</sup>, **Sarah J. Graves**<sup>2</sup>, **Dihong Gong**<sup>3</sup>, **Morteza Shahriari Nia**<sup>3</sup>, **Marion Le Bras**<sup>4</sup>, **Bonnie J. Dorr**<sup>4</sup>, **Peter Fontana**<sup>4</sup>, **Justin Gearhart**<sup>1</sup>, **Craig Greenberg**<sup>4</sup>, **Dave J. Harris**<sup>5</sup>, **Sugumar Arvind Kumar**<sup>3</sup>, **Agarwal Nishant**<sup>3</sup>, **Joshi Prarabdh**<sup>3</sup>, **Sundeep U. Rege**<sup>3</sup>, **Stephanie Ann Bohlman**<sup>Corresp., 2</sup>, **Ethan P. White**<sup>Corresp., 5</sup>, **Daisy Zhe Wang**<sup>Corresp., 3</sup>

<sup>1</sup> School of Natural Resources and Environment, University of Florida, Gainesville, Florida, United States

<sup>2</sup> School of Forest Resources and Conservation, University of Florida, Gainesville, Florida, United States

<sup>3</sup> Department of Computer and Information Science and Engineering, University of Florida, Gainesville, Florida, United States

<sup>4</sup> National Institute of Standards and Technology, Gaithersburg, United States

<sup>5</sup> Department of Wildlife Ecology and Conservation, University of Florida, Gainesville, Florida, United States

Corresponding Authors: Stephanie Ann Bohlman, Ethan P. White, Daisy Zhe Wang  
Email address: sbohlman@ufl.edu, ethanwhite@ufl.edu, daisyw@cise.ufl.edu

Ecology has reached the point where data science competitions, in which multiple groups solve the same problem using the same data by different methods, will be productive for advancing quantitative methods for tasks such as species identification from remote sensing images. We ran a competition to help improve three tasks that are central to converting images into information on individual trees: 1) crown segmentation, for identifying the location and size of individual trees; 2) alignment, to match ground truthed trees with remote sensing; and 3) species classification of individual trees. Six teams (composed of 16 individual participants) submitted predictions for one or more tasks. The crown segmentation task proved to be the most challenging, with the highest-performing algorithm yielding only 34% overlap between remotely sensed crowns and the ground truthed trees. However, most algorithms performed better on larger trees. For the alignment task, an algorithm based on minimizing the difference, in terms of both position and tree size, between ground truthed and remotely sensed crowns yielded a perfect alignment. In hindsight, this task was over simplified by only including targeted trees instead of all possible remotely sensed crowns. Several algorithms performed well for species classification, with the highest-performing algorithm correctly classifying 92% of individuals and performing well on both common and rare species. Comparisons of results across algorithms provided a number of insights for improving the overall accuracy in extracting ecological information from remote sensing. Our experience suggests that this kind of competition can benefit methods development in ecology and biology more broadly.

1 A data science challenge for converting airborne remote sensing data  
2 into ecological information

3 Sergio Marconi<sup>1</sup>, Sarah J. Graves<sup>2</sup>, Dihong Gong<sup>3</sup>, Morteza Shahriari Nia<sup>3</sup>, Marion Le Bras<sup>4</sup>,  
4 Bonnie J. Dorr<sup>4</sup>, Peter Fontana<sup>4</sup>, Justin Gearhart<sup>1</sup>, Craig Greenberg<sup>4</sup>, David J. Harris<sup>5</sup>, Arvind  
5 Kumar Sugumar<sup>3</sup>, Agarwal Nishant<sup>3</sup>, Prarabdh Joshi<sup>3</sup>, Sundeep Rege<sup>3</sup>, Stephanie Bohlman<sup>2\*</sup>,  
6 Ethan P. White<sup>5\*</sup>, Daisy Zhe Wang<sup>3\*</sup>

7

8 <sup>1</sup>School of Natural Resources and Environment, University of Florida, Gainesville, Florida,  
9 United States

10 <sup>2</sup>School of Forest Resources and Conservation, University of Florida, Gainesville, Florida,  
11 United States

12 <sup>3</sup>Department of Computer and Information Science and Engineering, University of Florida,  
13 Gainesville, Florida, United States

14 <sup>4</sup>National Institute of Standards and Technology, Gaithersburg, United States

15 <sup>5</sup>Department of Wildlife Ecology and Conservation, University of Florida, Gainesville, Florida,  
16 United States

17

18 \*Corresponding authors:

19 Stephanie Bohlman<sup>2</sup>, email: [sbohlman@ufl.edu](mailto:sbohlman@ufl.edu)

20 Ethan P. White<sup>5</sup>, email: [ethanwhite@ufl.edu](mailto:ethanwhite@ufl.edu)

21 Daisy Zhe Wang<sup>3</sup>, email: [daisyw@cise.ufl.edu](mailto:daisyw@cise.ufl.edu)

22

23

24 Abstract

25 Ecology has reached the point where data science competitions, in which multiple groups solve  
26 the same problem using the same data by different methods, will be productive for advancing  
27 quantitative methods for tasks such as species identification from remote sensing images. We ran  
28 a competition to help improve three tasks that are central to converting images into information  
29 on individual trees: 1) crown segmentation, for identifying the location and size of individual  
30 trees; 2) alignment, to match ground truth trees with remote sensing; and 3) species classification  
31 of individual trees. Six teams (composed of 16 individual participants) submitted predictions for  
32 one or more tasks. The crown segmentation task proved to be the most challenging, with the  
33 highest-performing algorithm yielding only 34% overlap between remotely sensed crowns and  
34 the ground truth trees. However, most algorithms performed better on larger trees. For the  
35 alignment task, an algorithm based on minimizing the difference, in terms of both position and  
36 tree size, between ground truth and remotely sensed crowns yielded a perfect alignment. In  
37 hindsight, this task was over simplified by only including targeted trees instead of all possible  
38 remotely sensed crowns. Several algorithms performed well for species classification, with the  
39 highest-performing algorithm correctly classifying 92% of individuals and performing well on  
40 both common and rare species. Comparisons of results across algorithms provided a number of  
41 insights for improving the overall accuracy in extracting ecological information from remote  
42 sensing. Our experience suggests that this kind of competition can benefit methods development  
43 in ecology and biology more broadly.

## 44 1. Introduction

45 In many areas of science and technology there are tasks for which solutions can be optimized  
46 using well-defined measures of success. For example, in the field of image analysis, the goal is  
47 to accurately characterize the largest proportion of images (Solomon & Breckon, 2010). When a  
48 clear measure of success can be defined, one approach to rapidly improving the methods used by  
49 the field is through open competitions (Carpenter, 2011). In these competitions, many different  
50 groups attempt to solve the same problem with the same data. This standardization of data and  
51 evaluation allows many different approaches to be assessed quickly and compared. Because the  
52 problems are well defined and data is cleaned and organized centrally, competitions can allow  
53 involvement by diverse participants, from those with domain expertise, to those in fields like  
54 modeling and machine learning.

55 In fields outside of ecology, these competitions have yielded rapid advances in the accuracy of  
56 many tasks. One well-known example of this is the ImageNET image classification competition  
57 (Krizhevsky et al., 2012). For the past five years, teams have competed in classifying 100,000s  
58 of images that has resulted in a major increase in classification accuracy from only 70% in 2010  
59 to 97% in 2017. This success has resulted in the rapid growth of competitions for solving  
60 common data science problems through both isolated competitions and major platforms like  
61 Kaggle (<https://www.kaggle.com/>). Kaggle has run over 200 competitions ranging from industry

62 challenges predicting sales prices of homes, to scientific questions like detecting lung cancer  
63 from lung scans. In general, life and environmental sciences, including ecology, have only  
64 recently begun to recognize the potential value of competitions. A few ecology-related  
65 competitions have been run recently, including competitions quantifying deforestation in the  
66 Amazon basin (<https://www.kaggle.com/c/planet-understanding-the-amazon-from-space>) and  
67 counting sea lions in Alaska ([https://www.kaggle.com/c/noaa-fisheries-steller-sea-lion-](https://www.kaggle.com/c/noaa-fisheries-steller-sea-lion-population-count)  
68 [population-count](https://www.kaggle.com/c/noaa-fisheries-steller-sea-lion-population-count)). However, these are far from common and, as a result, most ecologists are  
69 unaware of, and have had few opportunities to participate in, data science competitions.

70 In recent years, ecology has reached the point where these kinds of competitions could be  
71 productive. Large amounts of open data are increasingly available (Reichman et al. 2011,  
72 Hampton et al. 2013, Michener 2015) and areas of shared interest around which to center  
73 competitions are increasingly prominent. One of these shared areas of interest is converting  
74 remote sensing data into information on vegetation diversity, structure and function (Pettorelli et  
75 al. 2014, Pettorelli et al., 2017, Eddy et al., 2017). We ran a competition to improve three tasks  
76 that are central to converting airborne remote sensing (images and vertical structure  
77 measurements collected from airplanes) into the kinds of vegetation diversity and structure  
78 information traditionally collected by ecologists in the field: 1) crown segmentation, for  
79 identifying the location and size of individual trees (Zhen et al., 2016); 2) alignment to match  
80 ground truth data on trees with remote sensing data (Graves et al., in prep); and 3) species  
81 classification to identify trees to species (Fassnacht et al., 2016). If these three tasks can be  
82 conducted with a high degree of accuracy, it will allow scientists to map species locations over  
83 large areas, and use them to understand the factors governing the individual level behavior of  
84 natural systems at scales thousands of times larger than possible from traditional field work  
85 (Barbosa & Asner, 2017).

86 To create this competition, we used data from the National Ecological Observatory Network  
87 (NEON; Keller et al. 2008) funded by the U.S. National Science Foundation (NSF). NEON  
88 collects data from a wide range of ecological systems following standardized protocols. One of  
89 the core sets of observations comes from the Airborne Observation Platform (AOP) that collects  
90 high resolution LiDAR and hyperspectral images across ~10,000 ha for dozens of sites across the  
91 United States (<http://www.neonscience.org>). NEON also collects associated data on the  
92 vegetation structure at each site, which supports the building and testing of remote sensing based  
93 models. In addition to providing the openly available data needed for this competition, NEON  
94 also provides an ideal case for competitions because the methods are standardized across sites  
95 and data collection will be conducted at dozens of locations annually for the next 30 years. This  
96 means that the methodological improvements identified by the competition can be directly  
97 applied to hundreds of thousands of hectares of remotely sensed images and continual  
98 improvements can be made by regularly rerunning the competition. As a result, this competition  
99 has the potential to produce major gains in the quality of the ecological information that can be  
100 extracted from this massive data collection effort.

101 In addition to producing important improvements for NEON remote sensing products, this  
102 competition should also broadly benefit efforts to convert airborne remote sensing into  
103 ecological information. A major challenge in current assessments of airborne remote sensing  
104 tasks is determining whether published assessments of different methods generalize to the broad  
105 application of the methods as a whole, or are specific to the particular dataset and evaluation  
106 metrics being used. While this is a general problem for method comparison, it is particularly  
107 acute in many areas of remote sensing because: 1) most papers do not compare their methods to  
108 other approaches; 2) when comparisons are made it is typically between a new method and a  
109 single alternative; 3) different papers focus on different datasets; and 4) different papers often  
110 use different evaluation metrics and fail to specifically identify the best evaluation metric for a  
111 given task. Zhen et al. (2016) have highlighted the importance of changing this culture to  
112 produce extensive method comparisons using consistent data and evaluation metrics to drive the  
113 field of crown segmentation forward. By design, competitions provide single core datasets and  
114 consistent evaluation metrics to allow direct comparisons among many different approaches.

115 To capitalize on the benefits of competitions for overcoming barriers of comparing methods and  
116 determining how well different approaches to common data science task generalize, the National  
117 Institute of Standards and Technology (NIST) has been developing a Data Science Evaluation  
118 Series (DSE). This program has developed methodologies for evaluating progress in data science  
119 research through iterative examination of a range of problems, with the goal of devising a  
120 general evaluation paradigm to address data science problems that span diverse disciplines,  
121 domains, and tasks. As a part of the early stages of DSE, a pilot evaluation was run using traffic  
122 data, which was then followed by this competition on converting remote sensing data to  
123 information on trees. As a component of this endeavor, NIST researchers identified general  
124 classes of data science problems (Dorr et al., 2015; Dorr et al., 2016a, b) and produced a  
125 framework for evaluating methods both within an individual domain (like in this paper) and  
126 across domains (e.g., allowing algorithms for similar tasks to be applied to both traffic and  
127 ecological problems). This framework was used as the foundation for this competition including  
128 curating the datasets, developing the task and data descriptions, designing evaluation metrics,  
129 developing submission formats, and disseminating of participation information and rules.

130 Here we present the details of the initial run of this data science competition for converting  
131 remote sensing to data on individual trees. We present the details of the tasks and data, and  
132 summarize and synthesize the results from the participants. In a set of short accompanying  
133 papers and preprints, the participants describe the methods used, present detailed results for those  
134 methods, and discuss lessons learned and future directions for these methods (Anderson  
135 submitted, Dalponte et al. submitted, Taylor submitted, McMahon submitted, Sumison et al.  
136 submitted). Finally, we discuss the broad potential for competitions in ecology and the biological  
137 sciences more generally.

## 138 2. Materials & Methods

### 139 2.1. NEON data

140 We used NEON-AOP data (from year 2014) and field collected data (from years 2015-2017) for  
141 the Ordway-Swisher Biological Station (Domain D03, OSBS) in north-central Florida. The  
142 NEON field data was from 43 permanently established plots across the OSBS site, which are  
143 stratified across three land cover types (Homer et al. 2015). The field measurements were the  
144 NEON vegetation structure data that provides information on the stem location, taxonomic  
145 species, stem size, tree height, and in some cases two measurements of crown radius (Table 1).

146 Four NEON-AOP remote sensing data products were used; LiDAR point cloud data, LiDAR  
147 canopy height model (CHM), hyperspectral surface reflectance, and high resolution visible color  
148 (RGB) photographs (Table 1). The LiDAR point cloud data provide information about the  
149 vertical structure of the canopy. Data consists of a list of spatial 3D coordinates, with an average  
150 resolution of 4-6 points per square meter. The CHM data provides 1 m spatial resolution  
151 information on the spatial variation in canopy height. Hyperspectral data provides surface  
152 reflectance from 350-2500 nm at 1 m spatial resolution and allows development of spectral  
153 signatures to identify object categories. The RGB photographs provide 0.25 m spatial resolution  
154 information in the visible wavelengths. The higher spatial resolution relative to the other data  
155 products may be helpful to separate trees that are close to one another and to refine tree crown  
156 boundaries. The RGB data was the only data type not available for all plots (39 out of 43 total).  
157 NEON provides geographically registered files of these data products across the entire NEON  
158 site. The data was clipped to 80 x 80 m subsets to capture the full 40 x 40 m field plot with a 20  
159 m buffer on each side. The buffer was used to include any trees with their base in the plot but  
160 with a crown that fell outside of the NEON plot boundary.

### 161 2.2. Individual tree crown (ITC) field mapping data

162 Generating field-validated individual tree crowns (ITCs) required spatially matching individual  
163 trees measured in the field to the remote sensing image of their crowns taken from above the  
164 canopy. The ITCs were generated in the field on a tablet computer and GIS software. This  
165 process was done after NEON remote sensing and field data had been acquired and processed.  
166 First, the NEON images were loaded in a GIS application on a tablet computer that was  
167 connected to an external GPS device. The GIS software displayed the GPS location and the  
168 NEON digital images. Second, NEON plots were visited and field-technicians from our team  
169 located all tree crown that fell within a NEON plot and had branches that were in the upper  
170 canopy and visible in the NEON image. Third, with the aid of the GPS location, and the  
171 technicians' skills in visual image analysis, the crown boundaries of individual trees were  
172 digitized in the GIS application. While the LiDAR and RGB data was used to aid in tree crown  
173 delineation, the ITC polygons were made in reference to the hyperspectral data. This is important  
174 to consider when there is geographic misalignment among the 3 data products. The result of the

175 field mapping process was spatially explicit polygon objects that delineated the crown  
176 boundaries of individual trees. These polygons were linked to field data by the NEON  
177 identification number, or field-based species identification.

### 178 2.3. Train test split

179 Training data for the segmentation task consisted of a subset of 30 out of 43 plots (~70%). The  
180 ITCs were provided as ground truth to allow participants to apply supervised methods. Plots  
181 were selected to have a consistent 0.7 to 0.3 training-testing ratio both in number of plots, and  
182 number of ITCs (Table 2). The splitting resulted in a training dataset of 452 out of 613 ITCs.  
183 Since the OSBS NEON site is characterized by three different ecosystem types, we split the data  
184 accordingly to ensure each ecosystem was split in the 0.7 to 0.3 ratio. Separate polygon files  
185 were provided for each NEON plot. All ITC files had a variable number of polygons, and each  
186 polygon represented a single tree. LiDAR and hyperspectral derived data was made available to  
187 participants for all tasks. The RGB data were provided only when available. For the alignment  
188 task, we used only data from individual trees shared by the vegetation structure and the ITCs,  
189 resulting in a total of 130 entries. We split data in a 0.7 to 0.3 training-test ratio, following the  
190 same rationale described for segmentation. For the classification task we used data from all ITC  
191 crowns. Again, data were split in a 0.7 to 0.3 ratio. In this case, we stratified training-test  
192 samples by species labels (e.g. *Pinus palustris*, *Quercus laevis*). As a result, around 70% of the  
193 trees for each species belonged to the training set, the other 30% to the test set. We grouped  
194 species whose occurrences were less than 4 into a general category labelled as “Other”, because  
195 their individual numbers were considered too few to allow any learning. We consider the  
196 “Other” category potentially useful to discriminate rare, undefined species from the rest of the  
197 dataset.

### 198 2.4. Timeline and participants

199 The data science evaluation was announced one month in advance of making the data available  
200 (September 1, 2017), and participants were allowed to register until the final submission date  
201 (December 15, 2017). Participants could work on any or all of the tasks. There were two  
202 submission deadlines, with the first deadline providing an opportunity to get feedback on a  
203 submission evaluated on the test data before the final submission. A total of 84 groups showed  
204 interest in participating, 14 formally registered, and 6 teams submitted results. Teams came from  
205 a number of institutions including teams from outside the United States. The six teams were: 1)  
206 BYU, a team composed of 4 researchers from the Bioinformatics Research Group (BRG); 2)  
207 Conor, a team from University of Texas at Austin composed of a single researcher; 3) FEM, a  
208 team composed of 3 researchers of the Fondazione Edmund Mach (Italy); 4) GatorSense, a team  
209 composed of 5 members, all affiliated to University of Florida (but not involved in organizing  
210 the competition); 5) Shawn, a team composed of a single researcher at University of Florida; and  
211 6) StanfordCCB, a single researcher affiliated with Stanford University.



## 212 2.5. Competition Tasks

## 213 2.5.1. Segmentation

214 The crown segmentation task aims to determine the boundaries of tree crowns in an image.  
215 While image segmentation is a well developed field in computer science (Badrinarayanan et al.,  
216 2017, Saha & Panda, 2018), delineating tree crowns in a forest is a particularly complex task (Ke  
217 & Quackenbush, 2011; Bunting & Lucas, 2006). Most of the complexity is driven by the fact  
218 that individual crowns overlap, look similar to each other, and can show different shapes  
219 depending on the environment and developmental stage (Duncanson et al, 2014). The spatial  
220 resolutions of the NEON hyperspectral and LiDAR data (1m<sup>2</sup>) are also relatively low compared  
221 to crown sizes. In addition, these data are also different than most image data in that they have  
222 very high spectral resolution, which may facilitate the task of distinguishing neighboring tree  
223 crowns especially if coupled to LiDAR data. As a result of these complexities, there is no widely  
224 agreed upon solution to the crown segmentation problem, as widely described in Zhen et al.  
225 (2016). Different classes of algorithms perform best in different ecoregions, or even within a  
226 single forest. For example, the same method can perform well in an open canopy area and poorly  
227 in a closed canopy portion of the same stand.

228 For the segmentation task we asked participants to delineate tree crowns in the 80 x 80 m field-  
229 plot area using remote sensing data and the ITC polygons collected in the field (Figure 1). For a  
230 more detailed state of the art review, we point the reader to Zhen et al. (2016).

## 231 2.5.1.1. Performance metric

232 We used the mean pairwise Jaccard Coefficient,  $J(A,B)$ , as the performance metric for the  
233 segmentation task (Real & Vargas, 1996). The  $J(A,B)$  is a measure of similarity and diversity  
234 between pairs of objects, and is formulated as:

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|}.$$

235

236 Where A and B are respectively the observed and predicted ITCs. By definition, the  $J(A,B)$  is a  
237 value between 0 and 1, where 0 stands for no overlap, and 1 for a perfect match.

238 The score for the segmentation task is the average of the plot-level scores for each pair of  
239 crowns; that is, the average  $J(A,B)$  calculated on every measured ITC with the single most  
240 overlapping predicted crown. We used the Hungarian algorithm to match predicted and ground  
241 truth crowns. We chose this method because it is simple to interpret, does not require assignment  
242 of predicted crowns to specific ITCs by the participants, and provides a continuous measure. We  
243 penalized cases where predicted polygons overlapped with each other by disregarding the  
244 intersecting area in the numerator of the Jaccard Coefficient.

245 Although it was not an official scoring criterion, we also analyzed the confusion matrix of their  
246 predictions to detect how the errors were distributed. The confusion matrix is a table where  
247 predicted and ground truth labels are represented by columns and rows, respectively. In the  
248 context of crown delineation, labels are true positive, false positive, and true negative for each of  
249 the pixels. Given this information, we could determine and aggregate the number of false/true  
250 positives and negatives.

#### 251 2.5.1.2. Algorithms

252 Our baseline prediction consisted of applying the Chan-Vese algorithm (Chan et al., 2001) on the  
253 negative of the 1m<sup>2</sup> resolution canopy height model. Polygons boundaries were drawn by  
254 applying a segmentation mask to each predicted crown, and following their pixels' perimeter.  
255 Three groups participated in the segmentation task and each applied a different algorithm. The  
256 Conor group applied a three step method that first filtered pixels based on an greenness threshold  
257 (based on NDVI, the normalized difference vegetation index), then extracted local maxima from  
258 the canopy height model using a linear moving window, and finally ran a watershed  
259 segmentation seeded by the local maxima (McMahon et al. submitted). The FEM group applied a  
260 growing region algorithm based on relative distance and difference in reflectance between  
261 neighbor pixels (Dalponte et al., 2015, submitted). For this purpose, they used the hyperspectral  
262 images, and tuned the method by visual analysis on the training set. The Shawn group used a  
263 watershed algorithm on the CHM, filtering the scene by NDVI threshold (preprint).

#### 264 2.5.2. Alignment

265 Once crown location, position and shape are recognized, it is important to accurately identify  
266 which object in the images is linked to the data collected on the ground. Although both remote  
267 sensing and field data collection are georeferenced, these data products use different methods to  
268 acquire geolocation. Moreover, field data coordinates locate the central stem (trunk) position,  
269 instead of the crown's centroid, which can be offset from each other, especially in closed-canopy  
270 forests. The differences in stem and crown location could lead to substantial misalignment  
271 between the two products, and consequently to misattributed information that could affect the  
272 quality of further inference. This task is known as alignment and is the second step of the  
273 pipeline. The goal of alignment is to correctly label each tree crown polygon to a single tree in  
274 the ground data, thus allowing data collected on the ground (e.g., species identity, height, stem  
275 diameter, tree health) to be accurately associated with remote sensing data. For this round, we  
276 envisioned the alignment task as a 1:1 labelling problem (Figure 1). We provided ITC data for  
277 crowns sampled in the field only and asked participants to link each single ITC to a specific field  
278 label. We acknowledge that this is an oversimplification of the real problem because each single  
279 ground label could be potentially confused with several apparent crowns in proximity that were  
280 not included in the field-mapped ITC dataset.

### 281 2.5.2.1. Performance metric

282 Performance of matching field stem locations to ITCs was evaluated using the trace of the  
283 prediction matrix divided by the sum over the values in that matrix. This method was chosen  
284 based on the following reasoning. In the testing stage, suppose we have a set of remotely sensed  
285 data (ITC) denoted as  $\{p_n|n=1,\dots,N\}$ , and ground truth data denoted as  $\{g_n|n=1,\dots,N\}$ . We know in  
286 advance that there is a unique one-to-one mapping between the P and G sets. Without loss of  
287 generality, assume  $p_n$  should be mapped to  $g_n$  for  $n=1,\dots,N$ . For each data point  $p_i$ , a program  
288 predicts a non-negative confidence score that should be aligned with ground truth data point  $ij$ ,  
289 which forms a prediction matrix  $M = (m_{i,j})$  where  $i,j = 1,\dots,N$ . Then, the quality of prediction can  
290 be measured by the following scoring function:

$$291 \text{ score} = \frac{\text{trace}(M)}{\sum_{i,j} m_{i,j}}$$

292 where  $\text{trace}(\cdot)$  represents trace of a matrix and  $M$  represents the prediction matrix which has  
293 been aligned in the order which matches the ground truth.

### 294 2.5.2.2. Algorithms

295 Our baseline prediction was the application of naive Euclidean distance from the stem location to  
296 the centroid of the ITC. Two groups participated in this task and applied different algorithms.  
297 Both were based on the Euclidean distance between field stem and each of the ITCs included in  
298 the dataset. Euclidean distance was calculated by using East and North UTM spatial coordinates,  
299 as well as crown height and radius. The groups calculated these values using allometric  
300 relationships whenever tree height and crown size were missing from the field data. The Conor  
301 group used crown diameter as a measure of tree size (McMahon et al. submitted). Euclidean  
302 distances were adjusted for the average plot-level offset in the training data to compensate for  
303 location biases consistent within a plot. The FEM group applied the Euclidean distance based on  
304 spatial coordinates, tree height, and the crown radius as well (Dalponte et al. submitted). FEM  
305 used an allometric equation to estimate the crown radius from tree height. One of the main  
306 differences between the two methods was that FEM used a visual check on the results to  
307 manually correct points where the distance offset was too high.

### 308 2.5.3. Classification

309 A large number of ecological, environmental, and conservation-oriented questions depend on  
310 species identification. This includes efforts to conserve individual species, understand and  
311 maintain biodiversity, and incorporate the biosphere into global circulation models (Rocchini et  
312 al., 2015, Lees et al., 2018). Species identification is generally treated as a supervised problem,  
313 whose demand for labelled data is usually high. Linking remote sensing with field data would  
314 potentially provide species identification for thousands of trees, facilitating the building of a  
315 successful classifier. For this reason, we identified species classification as the last step of the

316 pipeline (Figure 1). Classifying trees species from remote sensing imagery is complicated by: (1)  
 317 highly unbalanced data; (2) features fundamental to differentiate among species that cannot be  
 318 perceived by the human eye; (3) contribution of the understory and soil to the image properties  
 319 for ITCs; and (4) data limitation, especially for rare species. A detailed description of the state of  
 320 the arts can be found in Fassnacht et al. (2016), and other methods borrowed by the field of  
 321 Image Vision in Wäldchen & Maler, (2017).

#### 322 2.5.3.1. Performance metric

323 We evaluated classification performance using two metrics. The first was rank-1 accuracy,  
 324 namely the fraction of crowns in the test set whose ground truth species identification  
 325 (species\_id) and genus identification (genus\_id) was assigned the highest probability by the  
 326 participant. It is calculated as:

$$rank1 = \frac{\sum_{n=1}^N \operatorname{argmax}_k(p(nk))}{N} == g_n$$

327

328 where  $g_n$  is the ground-truth class of crown  $i$ , and  $p_{nk}$  is the probability assigned by the  
 329 participant that crown  $i$  belongs to class  $k$ . This metric only considers whether the correct class  
 330 has the highest probability, not whether the probabilities are well-calibrated.

331 The second metric was the average categorical cross-entropy, defined as:

$$cost = \frac{-\sum_{n,k} \ln(p_{nk}) * \delta(g_n, k)}{N}$$

332

333 given that  $p_{nk} \neq 0$ , to avoid the singularity. The  $\delta(x, y)$  is an indicator function that takes value 1  
 334 when  $x = y$ . This metric rewards participants for submitting well-calibrated probabilities that  
 335 accurately reflect their uncertainty about which crowns belong to which class.

#### 336 2.5.3.2. Algorithms

337 Our baseline prediction was a classification based on probability distributions of species  
 338 frequency in the training data. The Conor group reduced the first 10 components of the  
 339 hyperspectral data and CHM information to three components, with two principal component  
 340 analysis (PCA) subsequently (McMahon et al. submitted). They applied a maximum likelihood  
 341 classifier to the test set to calculate the probability of each test tree to be a specific tree of the  
 342 training set. The class (species) of the tree in the test set was assigned by using the same label of  
 343 the individual tree with highest likelihood. The BRG group used a neural network multi-layer  
 344 perceptron on the hyperspectral images (Sumsion et al., submitted). Crown probabilities were  
 345 aggregated by averaging the pixel scale predicted probabilities. FEM applied a four step pipeline,  
 346 consisting of data normalization, Sequential Forward Floating feature selection, building of a

347 support vector machine classifier, and crown level aggregation by majority rule (Dalponte et al.  
348 submitted). The GatorSense group built a series of one-vs-one Applied Multiple Instance  
349 Adaptive Cosine Estimator (MI-ACE) classifiers (Zare et al., 2017) that automatically select the  
350 best subset of pixels to use for classification. Crown level probabilities were assigned by  
351 majority vote of pixel scale predictions. Finally, StanfordCCB group applied a six step pipeline  
352 (Anderson, submitted). Dimensionality reduction was performed using principal components  
353 analysis, and the first 100 components were retained. Pixels with high shade fractions were  
354 removed. Random Forest and Gradient Boosting multi-label classification algorithms were  
355 applied in a one-vs-all framework. Training species were under- or over-sampled to deal with  
356 label imbalance. Models' hyperparameters were determined using a grid search function, and  
357 prediction probabilities were calibrated using validation data. Finally, prediction probabilities  
358 were averaged between the two model ensembles.

### 359 3. Results

360 Overall, there was no single team that had a highest performing system across all three tasks. The  
361 FEM group achieved the highest evaluation scores for the segmentation and alignment tasks, but  
362 had a lower score for the classification task than the highest scoring group, StanfordCCB. In all  
363 three tasks, the highest scoring group scored substantially higher than the baseline. Given our  
364 evaluation data and metrics for each task, some groups performed better than the others.  
365 However, we may still be able to learn useful information or strategies from those teams that did  
366 not achieve the best performance on this specific competition configuration.

#### 367 3.1. Segmentation

368 This task had the lowest performance among the three tasks given our evaluation data and  
369 criteria (Figure 2). A segmentation that perfectly matched our field-delineated crowns would  
370 achieve of Jaccard score of 1.0000. All submissions performed well below the optimal score, but  
371 well above the baseline prediction. The highest-performing method, as determined by the Jaccard  
372 scoring function, achieved score of 0.3402 (Table 3). In comparison, our baseline system only  
373 has a score of 0.0863. All groups had more false positives compared to true positives, suggesting  
374 that all groups made polygons bigger than the field-based ITCs, on average (Figure 3). Only two  
375 groups, baseline and Connor (McMahon, submitted), had more false negatives than true positives  
376 indicating these approaches failed to segment some portion or all of a crown. Overall, the FEM  
377 group (Dalponte et al., submitted) had the best balance between minimizing false positive and  
378 negatives as well as the highest number of true positives, across trees with different crown size  
379 (Figure 4).

#### 380 3.2. Alignment

381 In this task, the FEM group again achieved the best performance, while the baseline system and  
382 the Conor group performed equally well. Surprisingly, the FEM group had the perfect accuracy

383 score of 1.0 (Figure 5). However, their pipeline is not fully automatable, and so may not be fully  
384 reproducible or scale to a significantly larger spatial extent. On the other hand, despite the  
385 similar structure to the automated part of FEM's method, Conor group did not perform any better  
386 than the baseline (Table 4).

### 387 3.3. Classification

388 We had the most participants in this task (6): BRG (Sumsion et al., submitted), Conor  
389 (McMahon, submitted), FEM (Dalponte et al., submitted), GatorSense, StanforCCB (Anderson,  
390 submitted) and our baseline system (Figure 6). For the evaluation criteria used in this  
391 competition, Cross Entropy loss (CE) and Rank-1 accuracy (Rank1), there was consistent  
392 ranking of all groups except our baseline system (Table 5). The top three groups in order were  
393 StanfordCCB, FEM, and Gatorsense (Figure 7). Conor and BRG outperformed our baseline  
394 system in Rank1 but not CE. Most of the difference in accuracy among groups was determined  
395 by ability in classifying species that were infrequent in the data set. In fact, all groups performed  
396 well in predicting the two most common species *Pinus palustris* (PIPA) and *Quercus laevis*  
397 (QULA), according to Rank1 scores (Figure 8). However, the three lowest-performing  
398 approaches (Baseline, BRG, and Conor) failed to predict all but these two species. StanfordCCB,  
399 FEM, and GatorSense were able to predict both PIPA and the rarest species (i.e. LIST and  
400 QUNI), but performed differently for the other species.

## 401 4. Discussion

402 The results of the competition are both promising and humbling, and the results for each task  
403 provide different lessons for how to improve both the conversion of remote sensing to ecological  
404 information, and the competition itself. An assessment of the results for each of the individual  
405 tasks is provided below.

### 406 4.1. Crown segmentation

407 The results of the crown segmentation task reveal the challenging nature of segmentation  
408 problems (Zhen et al., 2016). The highest-performing algorithms yielded only 34% overlap  
409 between the closest remotely sensed crowns and ground truth crowns mapped directly onto  
410 remote sensing imagery in the field. This suggests that crown segmentation algorithms have  
411 substantial room for improvement for precisely identifying individual crowns from remote  
412 sensing imagery.

413 By looking at the results across the three algorithms for this task, we can identify future  
414 directions for improvement. FEM, the best performing method, was the only method using  
415 hyperspectral data to perform segmentation, despite LiDAR data being used more commonly for  
416 segmentation (Zheng et al., 2016). This indicates that there is useful information in the  
417 hyperspectral data for classification. For example, the hyperspectral data may allow

418 distinguishing overlapping crowns from different species. As a result, some participants  
419 suggested that better segmentation may be achieved in the future by combining both  
420 hyperspectral and LiDAR derived information (McMahon submitted; Dalponte et al., submitted).  
421 However, it should be noted that the ground truth polygons were identified using the  
422 hyperspectral data (and not the LiDAR). This means that any misalignment resulting from  
423 preprocessing and orthorectification of the hyperspectral and LiDAR data would advantage  
424 hyperspectral data over LiDAR for this task.

425 This source of uncertainty is important beyond this competition because LiDAR data is typically  
426 used to perform segmentation, while hyperspectral data is usually used for classification. In case  
427 of misalignment, the exact segmentation on LiDAR would result in imperfect inclusion of  
428 hyperspectral pixels within associated crowns. As a result, LiDAR to hyperspectral misalignment  
429 should be taken into consideration when working with these data sources together and we will  
430 actively address it in future rounds of this Data Science Evaluation.

431 Exploring the accuracy of different segmentation algorithms more thoroughly reveals that  
432 uncertainty in delineating crowns is generally dependent on crown size (Figure 4). Crowns below  
433  $10 \text{ m}^2$  were poorly classified by all algorithms and most algorithms performed best for crown  
434 sizes over  $40 \text{ m}^2$ . This may be due to the fact that small crowns are often closer together, more  
435 heterogeneous in shape, and composed of fewer pixels. The highest-performing method, FEM's  
436 region growing algorithm, outperformed other algorithms on small and intermediate sized  
437 crowns. However, it performed worse than some other methods for the largest crowns. Conor's  
438 and Shawn's methods (preprint) generally performed best for larger crowns. This result shows  
439 the value of a comparative evaluation of different families of methods and suggests that creating  
440 ensembles of existing algorithms could result in better crown segmentation across the full range  
441 of tree sizes.

#### 442 4.2. Alignment

443 The results for the alignment tasks were promising. In fact, FEM's Euclidean distance based  
444 approach produced a perfect alignment between remotely sensed crowns and the stem location of  
445 individual trees. This precise match was accomplished by considering not only the position of  
446 the stem, but also the size of the crown. Adding the size of the crown was crucial for successful  
447 alignment because it allowed the algorithm to differentiate between multiple nearby stems based  
448 on differences in size. Using only Euclidean distance based on the position of the stems (the  
449 baseline) resulted in only a 48% alignment between stems and crowns. This perfect alignment is  
450 particularly encouraging because it used a statistical relationship between a standard field based  
451 measure of tree size (height) to estimate the size of the crown for the field data in cases where  
452 crown size was not measured. This means that the approach can be applied to all trees measured  
453 in the field, not just those where the less common direct measures of crown dimensions are  
454 performed. However, it is worth noting that FEM also performed a visual check of the  
455 alignments and shifted a few alignments manually based on this assessment (Dalponte et al.

456 submitted). This yielded meaningful improvements for crowns with misalignments of several  
457 meters or more (likely resulting from data entry or collection errors). While including manual  
458 steps is typically a concern for scaling up remote sensing predictions, it is less of an issue for  
459 alignment since this step is only important for model building, not prediction. That means that  
460 this step will typically only be applied to a few hundred or thousand trees making human  
461 involvement doable and potentially important.

462 While the alignment results are encouraging for linking remote sensing and ground truth data at  
463 the individual level, in hindsight, the degree of this success was also due in part to how we posed  
464 the problem for the competition. When selecting data for this task we only included trees that  
465 occurred in both the field and remote sensing data. In all cases, there were additional trees in the  
466 80 x 80 m image subsets that were not included in both the field and remote sensing data. This  
467 simplification resulted in overly sparse data compared to real-world situations where field data  
468 would need to be aligned against a full scene of remotely sensed crowns. Our original decisions  
469 made sense from an assessment perspective but failed to reflect the real-world complexity of the  
470 problem. We expect that including all trees in the scene will make the task more challenging. In  
471 the next round of the competition, we plan to include the remotely sensed crowns that lack  
472 corresponding field data to provide a clearer picture of the effectiveness in real-world situations.

#### 473 4.3. Classification

474 The species classification task was led by the StanfordCCB algorithm, which yielded the best  
475 overall performance with a categorical cross-entropy of 0.45 and a rank-1 accuracy of 92%  
476 (Figure 7). This is on the high end of classification accuracy rates reported for tree species  
477 identification from remote sensing (Fassnacht et al., 2016). This approach involved multiple  
478 preprocessing steps and an ensemble of Random Forest and Gradient Boosting multi-label  
479 classifications applied on each tree in a one-vs-all framework. A number of different models also  
480 performed well with rank-1 accuracies greater than 80% including Gatorsense, FEM, and Conor.  
481 StanfordCCB performed better in relation to other models when evaluated using categorical  
482 cross-entropy compared to rank-1 accuracy, which suggests that this method provides more  
483 accurate characterizations of uncertainty. Therefore, it is good at both identifying which species  
484 class a tree is most likely to belong to, and at knowing when it is unsure of which species to  
485 predict. This is a desirable property for a remote sensing model because good estimates of  
486 uncertainty allowing accurate error propagation into applications of those models. Exploring  
487 these results further by evaluating classifications for individual species (Figure 8; not part of the  
488 defined goals of the competition) shows that the StanfordCCB, FEM, and Gatorsense methods  
489 provide the best classifications for rare species, while other methods are only accurate for  
490 common species.

491 Interestingly, most of the groups that performed well developed multi-step methods that used  
492 data cleaning and dimensionality reduction. Outlier removal such as filtering dark or non-green  
493 pixels, seemed to be particularly important, likely because it allowed shadowed pixels or pixels



494 mixed with non-green vegetation like soil and wood, to be removed from the analysis. The  
495 Connor group averaged the spectra across all crown pixels and used structural information,  
496 namely crown radius and height range. Interestingly, averaging crown spectral information  
497 resulted in high predictability of the two most dominant classes, yet it was not a good strategy to  
498 predict rare species. This result suggests that clearing mixed noisy pixels may be particularly  
499 effective to better predict rare species. Likewise, adding structural features like crown radius  
500 may be useful in separating dominant classes. In general, the groups which performed best  
501 involved people with ecological expertise, which appeared useful in processing and selecting  
502 meaningful features from the data.

503 The other interesting aspect of the third task was the high participation. Five teams participated  
504 in this task compared to two teams for task 1 and three teams for task 2. We suspect that the  
505 higher level of participation was due to the task being the most straightforward, out-of-the-box,  
506 analysis. The relevant data was already extracted into a common tabular form meaning that most  
507 classification algorithms could be applied directly to the provided data. This makes the task  
508 easier for non-domain experts and suggests that standardizing tasks, so that a common set of  
509 algorithms can be readily applied to them, could result in greater participation in this type of  
510 competition and result in broad improvements across disciplines. This is the motivation behind a  
511 new NIST effort focused on algorithm transferability where the goal is to allow algorithms  
512 developed in one field to be applied to similar problems in other disciplines. The next iteration of  
513 the NIST DSE Series (Dorr et al., 2016a, b) will combine sets of related tasks from different  
514 domains to help drive this idea of algorithm transferability forward. Accomplishing this requires  
515 standardizing data formats to allow integration into a central automatic-scoring system. We are  
516 in the process of converting the data from this competition into schema provided by DARPA's  
517 Data-Driven Discovery of Models program (D3M) for this purpose.

518 Dealing with complex and non standard data types also highlights some of the challenges for  
519 data competitions in the environmental sciences. For example, most of the data in this  
520 competition is spatially explicit, a data type that does not completely generalize to more  
521 standard, non-spatial contexts, and involves file formats that many potential participants are not  
522 familiar with. We mitigated some of these challenges by cleaning and extracting simpler aspects  
523 of the data, but this also results in a loss of information relevant to the specific task. In fact, one  
524 participant found that the choices we had made to simplify the data limited their use of more  
525 advanced tools on the problem. In future rounds, we will seek to both provide simplified  
526 representations of the data that are accessible to many users and the full raw data that allow  
527 experts to employ tools appropriate to that data type.

#### 528 4.4. Insights from the competition

529 We developed and ran a data science competition on converting airborne remote sensing data  
530 into information on individual trees, with the goal of improving methods for using remote  
531 sensing to produce ecological information and accelerating methods development in ecology

532 more broadly. In developing this competition we took advantage of a major new source for open  
533 ecological and remote sensing data, the National Ecological Observatory Network (NEON).  
534 Because of the long-term large-scale nature of NEON's data collection, the results of the  
535 competition have the potential to go beyond general improvements in methods to yield  
536 immediate improvements in the quality of the ecological information that can be extracted from  
537 this massive data collection effort. The clearly defined goals, potential for general  
538 methodological improvements, and opportunity for immediate operationalization to produce data  
539 products that will be used by large numbers of scientists, makes this an ideal combination of  
540 problems and data for a data science competition.

541 We identified a single algorithm for each task that had the highest performance based on one or  
542 two performance criteria. While these algorithms showed the greatest promise for maximizing  
543 the evaluation criteria - for example providing the highest rank-1 classification accuracy for  
544 species identification - caution should be taken in focusing too much on a single method for  
545 several reasons. First, there are many different evaluation criteria that can be used depending on  
546 the specific application and ecological questions to be addressed. For example, in the evaluation  
547 criteria for the classification task, the correct identification of all trees was weighted equally,  
548 such that an algorithm that could correctly predict the common species would be favored over an  
549 algorithm that correctly predicted the rare species. Correct identification of the most common  
550 species may be the key goal for some ecological questions, such as producing maps of  
551 aboveground biomass. On the other hand, there may be other ecological questions for which  
552 equally good classification for all species is desirable. In this case, the training and test data may  
553 be chosen so that it is balanced among species, or weighting used in the evaluation criteria to  
554 increase the importance of identifying less common species (Graves et al., 2016; Anderson  
555 submitted did this for this competition). For some biodiversity assessments, the optimization for  
556 the species classification task may be more focused on identifying rare species, a single exotic  
557 species, or identifying species that are outliers, and potentially "new" or unusual species in the  
558 system (Baldeck et al. 2015). The evaluation criteria for these alternative goals would differ from  
559 the ones used in this competition.

560 In addition to performing differently for a variety of specific tasks, different algorithms may vary  
561 in applicability and performance in different ecosystems or when using different types of field  
562 data. This competition used individual tree crowns from forest ecosystems, which are multi-pixel  
563 objects, as the unit of observation. However, other ecosystems, such as grasslands, prairies, open  
564 savannas, and shrublands, are dominated by plant species whose size is below the resolution of  
565 an individual pixel. NEON provides extensive data sets on the presence and cover of small plant  
566 species that if linked with NEON AOP data, could be used to generate landscape maps of these  
567 species. At a number of sites these sub-pixel plant species are the dominant plants at the site.  
568 Working across all NEON sites will therefore require algorithms that can perform alignment and  
569 identification of both super- and sub-pixel resolutions, a complex task that may change which  
570 algorithms perform best. For example, one of the approaches used in this competition,

571 GatorSense's multiple instance classification, had slightly lower performance than the highest-  
572 performing method in the species classification task, but has the flexibility to be used for the  
573 alignment and subpixel detection of small plant species presence and cover (Zare et al., 2017).  
574 This suggests that despite not being the highest-performing method in the competition, it is a  
575 promising route forward for the more general task.

576 For competitions like this one to be most effective in facilitating rapid methodological  
577 improvement of a field, it is important that the details of each teams' analysis be described in  
578 detail and easy to reproduce. This allows for researchers to quickly integrate the advances made  
579 by other participants into their own workflows. We accomplished this for this competition in  
580 three ways. First, all of the data is openly available under an open license (ECODSE group  
581 2017). Second, all authors wrote short papers describing the detailed methods employed in their  
582 analyses and these papers are published as part of collection associated with this paper (link to  
583 PeerJ collection). Finally, all authors posted their code openly on GitHub and linked it in their  
584 contributions. One author (Anderson 2018) even encouraged other researchers to use and further  
585 improve on their method with the hope of collaboratively improving the use of remote sensing  
586 for species classification. Having access to a growing number of fully reproducible open  
587 pipelines evaluated on the same data will be a powerful instrument improving the methods used  
588 in converting remote sensing into ecological information.

589 We plan to continue to run this competition, updating the specifics of the tasks to help advance  
590 the science of converting remote sensing to information on individual trees. In the next iteration  
591 of this competition, we plan to address the fact that remote sensing models for identification of  
592 species and other key ecosystem traits are usually developed at individual sites (Zhen et al.,  
593 2016, Fassnacht et al., 2016), which tend to make them site-specific and leads to a profusion of  
594 locally optimized methods that do not transfer well to other locations. For standardized data  
595 collection efforts like NEON, algorithms and models that perform well across sites are critical.  
596 To facilitate advances in this area we will include data from multiple NEON sites in future  
597 competitions with the goals of developing algorithms with high cross-site performance and  
598 comparing the performance of cross-site and site-specific algorithms.

## 599 5. Conclusions

600 The results of this competition are encouraging both for the specific scientific tasks involved and  
601 for the use of competitions in ecology and science more broadly. The highest performing  
602 algorithms are indicative of the potential for using remote sensing models to obtain reasonable  
603 estimates of the location and species identity of individual trees. The competition results help  
604 highlight the components of this process that have good existing solutions as well as those most  
605 in need of improvement. Promising areas for future development include the ensemble of crown  
606 segmentation algorithms that perform well for small vs. large crowns. In cases with clearly  
607 defined outcomes, science would benefit from the increased use of competitions as a way to  
608 quickly determine and improve on the highest-performing methods currently available.

## 609 6. Acknowledgements

610 The National Ecological Observatory Network is a program sponsored by the National Science  
611 Foundation and operated under cooperative agreement by Battelle Memorial Institute. This  
612 material is based in part upon work supported by the National Science Foundation through the  
613 NEON Program.

614 These results are not to be construed or represented as endorsements of any participants system,  
615 methods, or commercial product, or as official findings on the part of NIST or the U.S.  
616 Government. Certain commercial equipment, instruments, software, or materials are identified in  
617 this paper in order to specify the experimental procedure adequately. Such identification is not  
618 intended to imply recommendation or endorsement by NIST, nor is it intended to imply that the  
619 equipment, instruments, software or materials are necessarily the best available for the purpose.

## 620 7. References

621 Badrinarayanan, V., Kendall, A. and Cipolla, R., 2017. Segnet: A deep convolutional encoder-  
622 decoder architecture for image segmentation. *IEEE transactions on pattern analysis and machine*  
623 *intelligence*, 39(12), pp.2481-2495.

624 Baldeck, C. A., Asner, G. P., Martin, R. E., Anderson, C. B., Knapp, D. E., Kellner, J. R., &  
625 Wright, S. J. (2015). Operational tree species mapping in a diverse tropical forest with airborne  
626 imaging spectroscopy. *PLoS One*, 10(7), e0118403.

627 Barbosa, J.M. and Asner, G.P., 2017. Prioritizing landscapes for restoration based on spatial  
628 patterns of ecosystem controls and plant-plant interactions. *Journal of Applied Ecology*, 54(5),  
629 pp.1459-1468.

630 Bunting, P. and Lucas, R., 2006. The delineation of tree crowns in Australian mixed species  
631 forests using hyperspectral Compact Airborne Spectrographic Imager (CASI) data. *Remote*  
632 *Sensing of Environment*, 101(2), pp.230-248.

633 Carpenter, J., 2011. May the best analyst win. *Science (New York, NY)*, 331(6018), pp.698-699.

634 Chan, Tony F., and Luminita A. Vese. "Active contours without edges." *IEEE Transactions on*  
635 *image processing* 10.2 (2001): 266-277.

636 Dalponte, M., Reyes, F., Kandare, K. and Gianelle, D., 2015. Delineation of individual tree  
637 crowns from ALS and hyperspectral data: a comparison among four methods. *European Journal*  
638 *of Remote Sensing*, 48(1), pp.365-382.

639 Dorr, Bonnie J., Craig S. Greenberg, Peter Fontana, Mark Przybocki, Marion Le Bras, Cathryn  
640 Ploehn, Oleg Aulov, Martial Michel, E. Jim Golden, and Wo Chang, "The NIST IAD Data

- 641 Science Research Program," Proceedings of the IEEE International Conference on Data Science  
642 and Advanced Analytics, Paris, France, pp. 1-10, 2015.
- 643 Dorr, Bonnie J., Craig S. Greenberg, Peter Fontana, Mark Przybocki, Marion Le Bras, Cathryn  
644 Ploehn, Oleg Aulov, Martial Michel, E. Jim Golden, Wo Chang, "A New Data Science Research  
645 Program: Evaluation, Metrology, Standards, and Community Outreach," International Journal of  
646 Data Science and Analytics, Volume 1, Issue 3, Springer, pp. 1-21, October 2016b.
- 647 Dorr, Bonnie J., Peter C. Fontana, Craig S. Greenberg, Marion Le Bras, Mark Przybocki,  
648 "Evaluation-Driven Research in Data Science: Leveraging Cross-Field Methodologies", In  
649 Proceedings of the IEEE International Conference on Big Data (IEEE BigData), Washington,  
650 DC, December, 2016a.
- 651 Duncanson, L.I., Cook, B.D., Hurtt, G.C. and Dubayah, R.O., 2014. An efficient, multi-layered  
652 crown delineation algorithm for mapping individual tree structure across multiple ecosystems.  
653 Remote Sensing of Environment, 154, pp.378-386.
- 654 ECODSE group. 2017. ECODSE competition training set [Data set]. Zenodo.  
655 <http://doi.org/10.5281/zenodo.1206101>
- 656 Eddy, I.M., Gergel, S.E., Coops, N.C., Henebry, G.M., Levine, J., Zerriffi, H. and Shibkov, E.,  
657 2017. Integrating remote sensing and local ecological knowledge to monitor rangeland  
658 dynamics. Ecological indicators, 82, pp.106-116.
- 659 Fassnacht, F.E., Latifi, H., Stereńczak, K., Modzelewska, A., Lefsky, M., Waser, L.T., Straub, C.  
660 and Ghosh, A., 2016. Review of studies on tree species classification from remotely sensed data.  
661 Remote Sensing of Environment, 186, pp.64-87.
- 662 Gatzliolis, D., Fried, J.S. and Monleon, V.S., 2010. Challenges to estimating tree height via  
663 LiDAR in closed-canopy forests: a parable from western Oregon. Forest Science, 56(2), pp.139-  
664 155.
- 665 Graves, S.J., Asner, G.P., Martin, R.E., Anderson, C.B., Colgan, M.S., Kalantari, L. and  
666 Bohlman, S.A., 2016. Tree species abundance predictions in a tropical agricultural landscape  
667 with a supervised classification model and imbalanced data. Remote Sensing, 8(2), p.161.
- 668 Hampton, S.E., Strasser, C.A., Tewksbury, J.J., Gram, W.K., Budden, A.E., Batcheller, A.L.,  
669 Duke, C.S. and Porter, J.H., 2013. Big data and the future of ecology. Frontiers in Ecology and  
670 the Environment, 11(3), pp.156-162.
- 671 Homer, C.G., Dewitz, J.A., Yang, L., Jin, S., Danielson, P., Xian, G., Coulston, J., Herold, N.D.,  
672 Wickham, J.D., and Megown, K., 2015, Completion of the 2011 National Land Cover Database  
673 for the conterminous United States-Representing a decade of land cover change information.  
674 Photogrammetric Engineering and Remote Sensing, v. 81, no. 5, p. 345-354.

- 675 Ke, Y. and Quackenbush, L.J., 2011. A review of methods for automatic individual tree-crown  
676 detection and delineation from passive remote sensing. *International Journal of Remote Sensing*,  
677 32(17), pp.4725-4747.
- 678 Keller, M., Schimel, D.S., Hargrove, W.W. and Hoffman, F.M., 2008. A continental strategy for  
679 the National Ecological Observatory Network. *Frontiers in Ecology and the Environment*, 6(5),  
680 pp.282-284.
- 681 Krizhevsky, A., Sutskever, I. and Hinton, G.E., 2012. Imagenet classification with deep  
682 convolutional neural networks. In *Advances in neural information processing systems* (pp. 1097-  
683 1105).
- 684 Lees, K.J., Quaife, T., Artz, R.R.E., Khomik, M. and Clark, J.M., 2018. Potential for using  
685 remote sensing to estimate carbon fluxes across northern peatlands—A review. *Science of the*  
686 *Total Environment*, 615, pp.857-874.
- 687 Michener, William K. "Ecological data sharing." *Ecological Informatics* 29 (2015): 33-44.
- 688 NIST SRE. "The 2013-2014 Speaker Recognition i-vector Machine Learning Challenge."  
689 [https://www.nist.gov/sites/default/files/documents/itl/iad/mig/sre-ivectorchallenge\\_2013-11-](https://www.nist.gov/sites/default/files/documents/itl/iad/mig/sre-ivectorchallenge_2013-11-18_r0.pdf)  
690 [18\\_r0.pdf](https://www.nist.gov/sites/default/files/documents/itl/iad/mig/sre-ivectorchallenge_2013-11-18_r0.pdf). 2013.
- 691 National Ecological Observatory Network. 2014. Data Products NEON.DP1.10098,  
692 NEON.DP1.30010, NEON.DP3.30015, NEON.DP1.30008. Provisional data downloaded from  
693 <http://data.neonscience.org> on 26 Jan 2016. Battelle, Boulder, CO, USA
- 694 Pettorelli, N., H. Nagendra, R. Williams, D. Rocchini, and E. Fleishman. 2015. A new platform  
695 to support research at the interface of remote sensing, ecology and conservation. *Remote Sens.*  
696 *Ecol. Conserv.* 1, 1–3.
- 697 Pettorelli, N., Nagendra, H., Rocchini, D., Rowcliffe, M., Williams, R., Ahumada, J., De Angelo,  
698 C., Atzberger, C., Boyd, D., Buchanan, G. and Chauvenet, A., 2017. Remote Sensing in Ecology  
699 and Conservation: three years on. *Remote Sensing in Ecology and Conservation*, 3(2), pp.53-56.
- 700 Prange, John D. "Evaluation Driven Research: The Foundation of the TIPSTER Text Program,"  
701 in *Proceedings of the Workshop on TIPSTER '96*, Association for Computational Linguistics,  
702 pp. 13-22, 1996.
- 703 Real, R. and Vargas, J.M., 1996. The probabilistic basis of Jaccard's index of similarity.  
704 *Systematic biology*, 45(3), pp.380-385.
- 705 Reichman, O. James, Matthew B. Jones, and Mark P. Schildhauer. "Challenges and opportunities  
706 of open data in ecology." *Science* 331, no. 6018 (2011): 703-705.

- 707 Rocchini, D., Andreo, V., Förster, M., Garzon-Lopez, C.X., Gutierrez, A.P., Gillespie, T.W.,  
708 Hauffe, H.C., He, K.S., Kleinschmit, B., Mairota, P. and Marcantonio, M., 2015. Potential of  
709 remote sensing to predict species invasions: a modelling perspective. *Progress in Physical*  
710 *Geography*, 39(3), pp.283-309.
- 711 Saha, M. and Panda, C., 2018. A Review on Various Image Segmentation Techniques for Brain  
712 Tumor Detection.
- 713 Solomon, C. and Breckon, T., 2011. *Fundamentals of Digital Image Processing: A practical*  
714 *approach with examples in Matlab*. John Wiley & Sons.
- 715 Wäldchen, J. and Mäder, P., 2018. Plant species identification using computer vision techniques:  
716 A systematic literature review. *Archives of Computational Methods in Engineering*, 25(2),  
717 pp.507-543.
- 718 Zare, A., Jiao, C. and Glenn, T., 2017. Discriminative multiple instance hyperspectral target  
719 characterization. *IEEE transactions on pattern analysis and machine intelligence*.
- 720 Zhen, Z., Quackenbush, L.J. and Zhang, L., 2016. Trends in automatic individual tree crown  
721 detection and delineation—evolution of lidar data. *Remote Sensing*, 8(4), p.333.

**Table 1** (on next page)

*Data products and sources (National Ecological Observatory Network, 2016).*

*Information about data products can be found on the NEON data products catalogue (<http://data.neonscience.org/data-product-catalog> ).*



Name	NEON data product ID	Data date	How it was used
<a href="#">Woody plant vegetation structure</a>	<a href="#">NEON.DP1.10098</a>	2015	Task 2 vegetation structure
<a href="#">Spectrometer orthorectified surface directional reflectance - flightline</a>	<a href="#">NEON.DP1.30008</a>	2014	Task 1, 2, and 3 RS data (Hyperspectral)
<a href="#">Ecosystem structure</a>	<a href="#">NEON.DP3.30015</a>	2014	Task 1, 2, and 3 RS data (Canopy height model)
<a href="#">High-resolution orthorectified camera imagery</a>	<a href="#">NEON.DP1.30010</a>	2014	Task 1, 2, and 3 RS data (RGB photos)
Field ITC	Internal	2017	Task 1 ITC data; Task 3 to extract pixels per each crown

1

**Table 2** (on next page)

*Overview of train-test data split by task and ecosystem type.*

*The columns present respectively the number of NEON plots (Plots) and Individual Tree Crowns (ITC) provided per task and ecosystem type. EF, Evergreen Forest; EHW, Emergent Herbaceous Wetland; WWET, Woody Wetland.*

	Task 1			Task 2			Task 3	
	Plots	ITC		Plots	ITC		Plots	ITC
	Train							
EF	22	349		17	82		22	349
EHW	2	52		0	0		2	52
WWET	6	9		1	2		6	9
Total	30	452		19	84		30	452
	Test							
EF	9	144		7	28		9	144
EHW	1	21		0	0		1	21
WWET	3	7		1	2		3	7
Total	13	172		8	30		13	172

1

**Table 3** (on next page)

*Comparison of Jaccard scores among submissions and baseline*

Task 1: Crown Delineation		
Rank	Participant	Score
#1	FEM	0.3402
#2	Conor	0.184
#3	Shawn	0.0555
	Baseline	0.0863

1

**Table 4** (on next page)

*Comparison of alignment accuracy among submissions and baseline.*

Task 2: Crown Alignment		
Rank	Participant	Score
#1	FEM	1
#2	Conor	0.48
	Baseline	0.48

1

**Table 5** (on next page)

*Comparison of classification performance on categorical cross-entropy and rank-1 accuracy among submissions and baseline.*



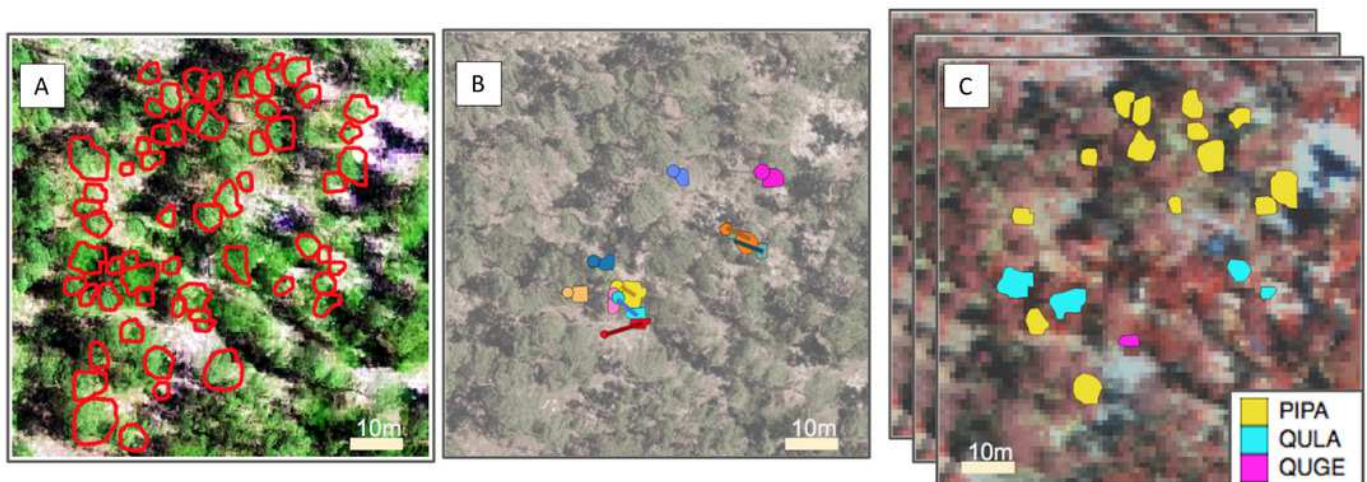
Task 3: Species Classification			
Rank	Participant	Score (Cross Entropy)	Score (Rank-1 Accuracy)
#1	StanfordCCB	0.4465	0.9194
#2	FEM	0.8769	0.88
#3	GatorSense	0.9386	0.864
#4	Conor	1.2247	0.8226
#5	BRG	1.4478	0.688
	Baseline	1.1306	0.6667

1

# Figure 1

*Representation of the pipeline for the three competition tasks.*

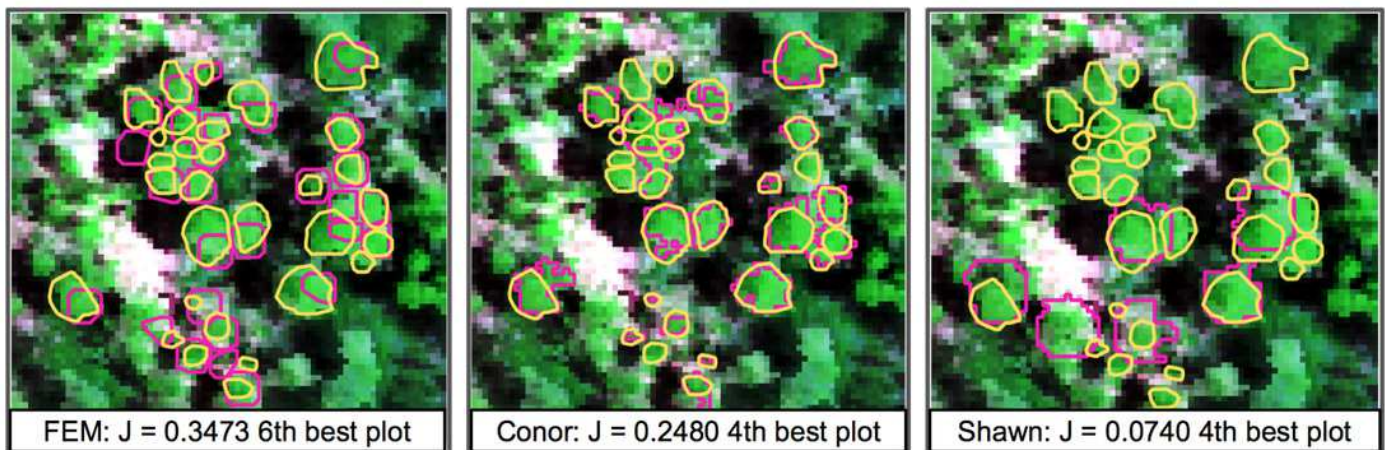
*From left to right, (A) Segmentation shows field ITC (red); the background is a composite of the hyperspectral data overlaid by LiDAR CHM. (B) Alignment shows stem locations scaled by stem diameter (circles) and field ITCs (irregular polygons) overlaid over a desaturated RGB image. Both ITCs and stem locations colored by stem identity. Lines indicating the offset between crowns and stems. (C) Classification shows field ITCs colored by species code. The background is a false-color composite of the hyperspectral data.*



## Figure 2

Sample of the participants' algorithm performance on average on plot 41, ranked around the median highest in performance for all the 3 groups (ranking 6th, 4th, and 4th respectively).

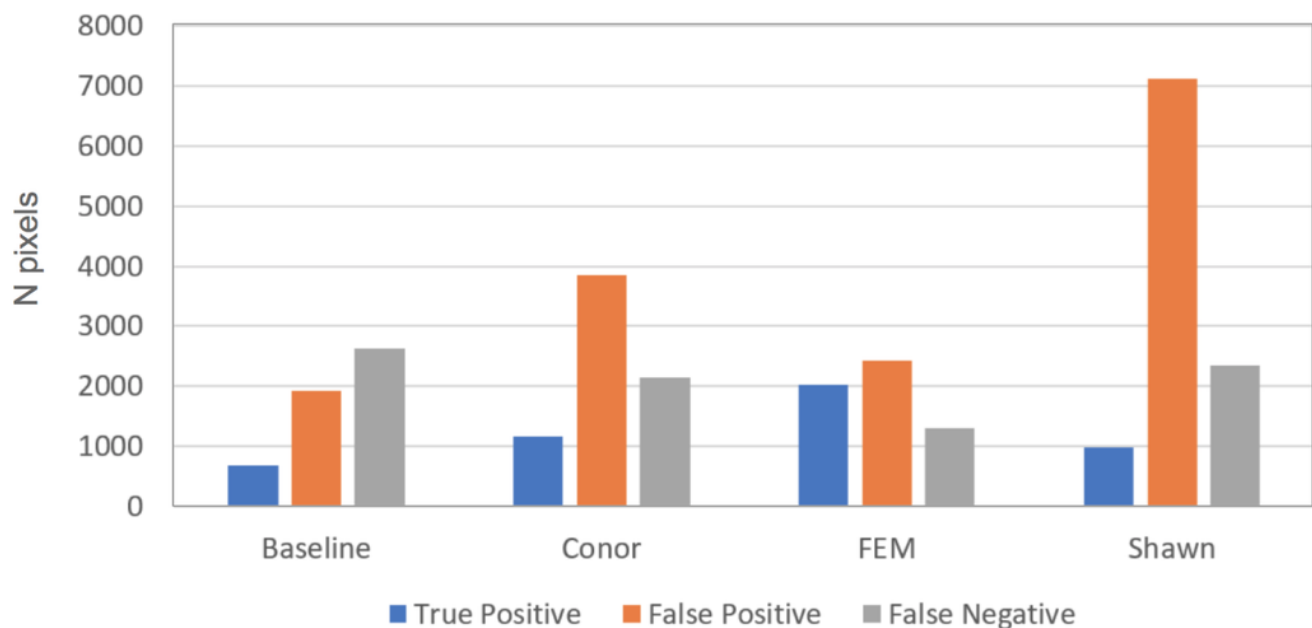
Yellow polygons represent ground truth ITCs, magenta the predicted ITCs. The background image is a composite of the hyperspectral data overlaid by LiDAR CHM.



## Figure 3

*Summary of error types for the crown segmentation task, using the 2 by 2 confusion matrix.*

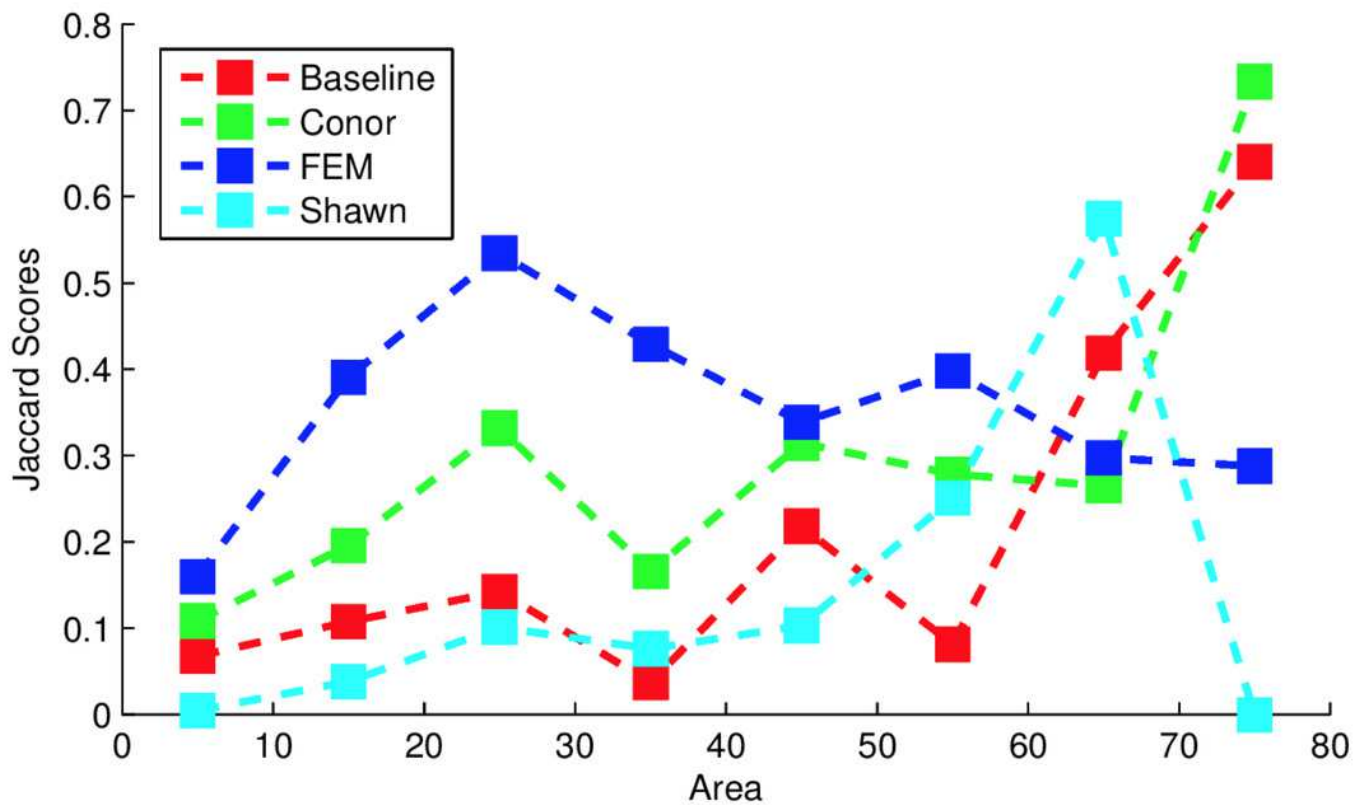
*Although presented in this figure, in the current competition evaluation criteria, we did not use false negatives, since the ground truth ITCs did not cover the entire image area. For this reason, the number of pixels obtained by summing the three columns per each group do not necessarily match among submissions.*



## Figure 4

*Jaccard score for crown segmentation as a function of the size (area) of the tree crown.*

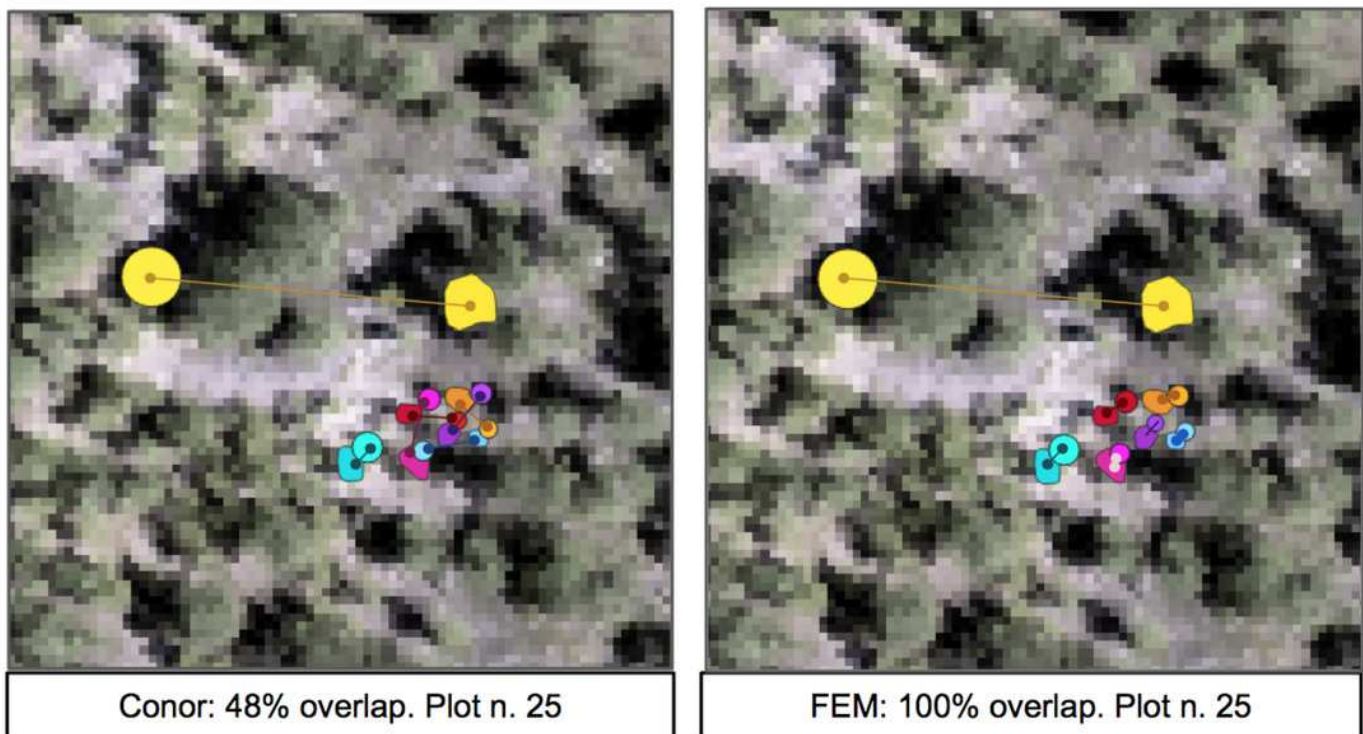
*Jaccard scores for individual trees are binned into size classes and averaged.*



## Figure 5

*Sample of the participants' algorithm performance on plot 25, for the two competing groups.*

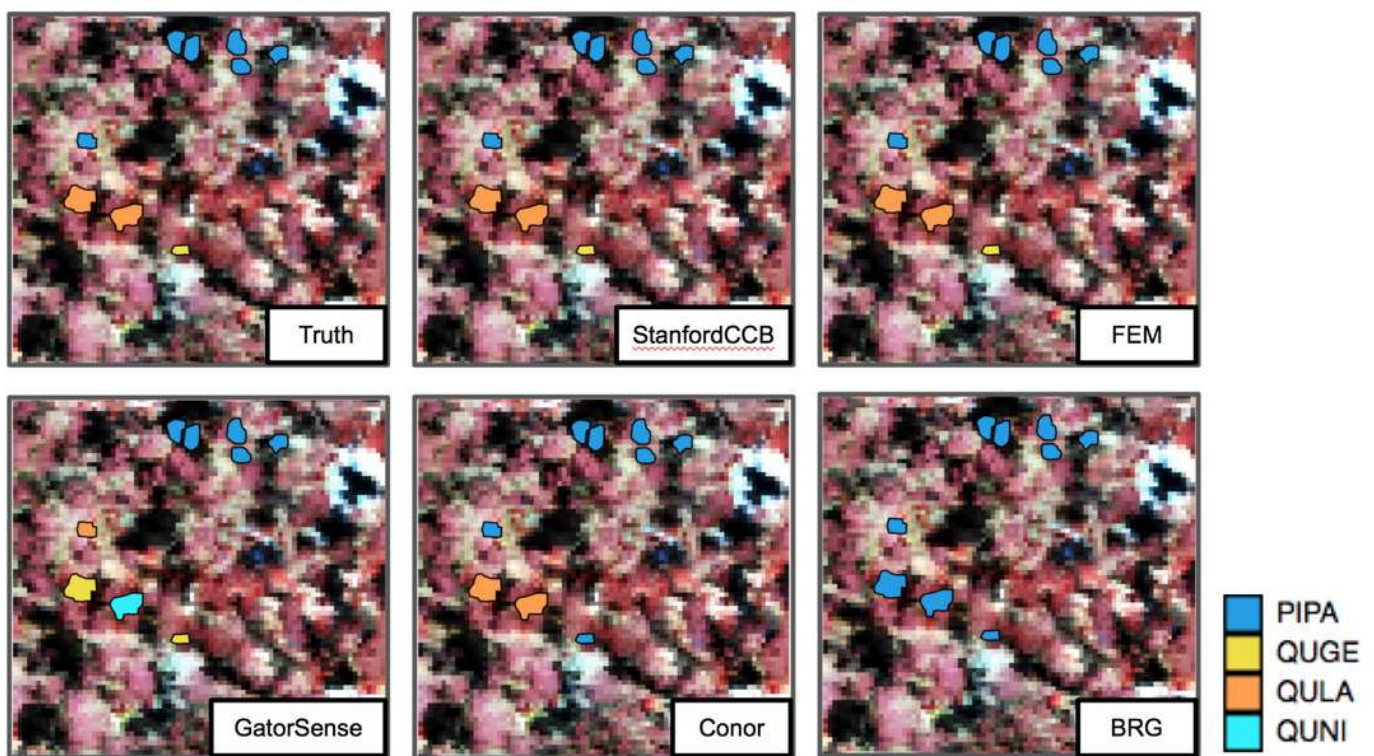
*Plot 25 was chosen for visualization because of the presence of one crown highly misaligned. Data shown are stem locations (circles) scaled by diameter at breast height; field ITCs (polygons); euclidean distances between the two data sources with same stem identity (solid line). ITCs, stem, and distances colored by stem identity. Images background is desaturated hyperspectral composite image.*



## Figure 6

*Performance of species classification in a plot that is relatively diverse in species composition.*

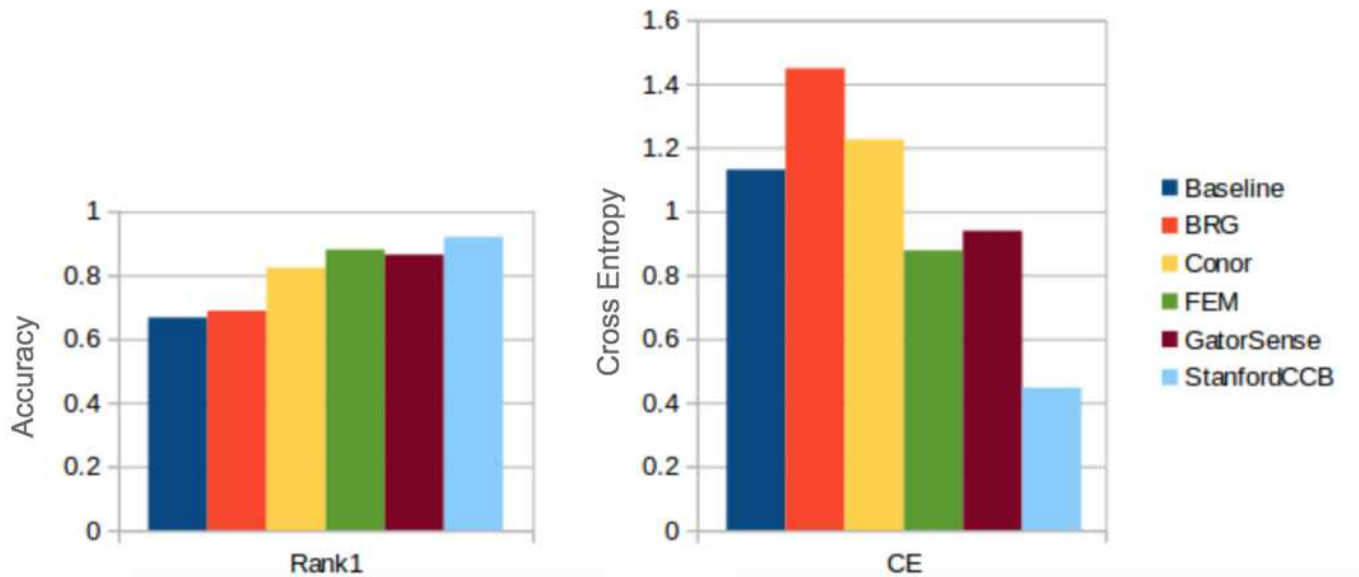
*Field ITCs are colored by species code. The background is a false-color composite of the hyperspectral data.*



# Figure 7

*Classification performance comparison.*

(A) Rank 1 accuracy; (B) Categorical cross-entropy.





## Figure 8

Comparison of Rank-1 classification accuracy by species.

The number in square bracket is number of training samples.

