



UNIVERSITY OF LEEDS

This is a repository copy of *Comparing the performance of flat and hierarchical Habitat/Land-Cover classification models in a NATURA 2000 site*.

White Rose Research Online URL for this paper:
<http://eprints.whiterose.ac.uk/125666/>

Version: Accepted Version

Article:

Gavish, Y orcid.org/0000-0002-6025-5668, O'Connell, J, Marsh, CJ et al. (4 more authors) (2018) Comparing the performance of flat and hierarchical Habitat/Land-Cover classification models in a NATURA 2000 site. *ISPRS Journal of Photogrammetry and Remote Sensing*, 136. pp. 1-12. ISSN 0924-2716

<https://doi.org/10.1016/j.isprsjprs.2017.12.002>

© 2017, Elsevier. This manuscript version is made available under the CC-BY-NC-ND 4.0 license <http://creativecommons.org/licenses/by-nc-nd/4.0/>

Reuse

This article is distributed under the terms of the Creative Commons Attribution-NonCommercial-NoDerivs (CC BY-NC-ND) licence. This licence only allows you to download this work and share it with others as long as you credit the authors, but you can't change the article in any way or use it commercially. More information and the full terms of the licence here: <https://creativecommons.org/licenses/>

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.



eprints@whiterose.ac.uk
<https://eprints.whiterose.ac.uk/>

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

TITLE:

Comparing the performance of flat and hierarchical habitat/land-cover classification models in a NATURA 2000 site.

AUTHORS:

Yoni Gavish^{a,*}	gavishyoni@gmail.com
Jerome O'Connell^b	jerome.oconnell@ucd.ie
Charles J. Marsh^a	charliem2003@gmail.com
Cristina Tarantino^c	cristina.tarantino@iia.cnr.it
Palma Blonda^c	palma.blonda@iia.cnr.it
Valeria Tomaselli^d	valeria.tomaselli@ibbr.cnr.it
William E. Kunin^a	W.E.Kunin@leeds.ac.uk

AFFILIATIONS:

^a School of biology, University of Leeds, Leeds, LS2 9JT, United Kingdom
^b School of Biosystems and Food Engineering, University College Dublin, D04 N2E5, Ireland
^c Institute of Atmospheric Pollution Research (IIA), National Research Council (CNR), c/o Interateneo Physics Department, Via Amendola 173, 70126 Bari, Italy
^d Institute of Biosciences and BioResources (IBBR), National Research Council (CNR-IBBR), via G.Amendola 165/A 70126, Bari, Italy

CORRESPONDENCE:

* Corresponding author: Yoni Gavish, School of biology, University of Leeds, Leeds, LS2 9JT, United Kingdom. Mobile: 00-44-75-9999-1988. Email: gavishyoni@gmail.com

23 **Abstract**

24 The increasing need for high quality Habitat/Land-Cover (H/LC) maps has triggered considerable
25 research into novel machine-learning based classification models. In many cases, H/LC classes follow
26 pre-defined hierarchical classification schemes (e.g., CORINE), in which fine H/LC categories are
27 thematically nested within more general categories. However, none of the existing machine-learning
28 algorithms account for this pre-defined hierarchical structure. Here we introduce a novel Random
29 Forest (RF) based application of hierarchical classification, which fits a separate local classification
30 model in every branching point of the thematic tree, and then integrates all the different local
31 models to a single global prediction. We applied the hierarchal RF approach in a NATURA 2000 site in
32 Italy, using two land-cover (CORINE, FAO-LCCS) and one habitat classification scheme (EUNIS) that
33 differ from one another in the shape of the class hierarchy. For all 3 classification schemes, both the
34 hierarchical model and a flat model alternative provided accurate predictions, with kappa values
35 mostly above 0.9 (despite using only 2.2-3.2% of the study area as training cells). The flat approach
36 slightly outperformed the hierarchical models when the hierarchy was relatively simple, while the
37 hierarchical model worked better under more complex thematic hierarchies. Most misclassifications
38 came from habitat pairs that are thematically distant yet spectrally similar. In 2 out of 3 classification
39 schemes, the additional constraints of the hierarchical model resulted with fewer such serious
40 misclassifications relative to the flat model. The hierarchical model also provided valuable
41 information on variable importance which can shed light into “black-box” based machine learning
42 algorithms like RF. We suggest various ways by which hierarchical classification models can increase
43 the accuracy and interpretability of H/LC classification maps.

44

45 **Key Words**

46 Classification, machine-learning, hierarchical models, random forest, NATURA 2000, Habitat/Land-
47 Cover

48

49

50

51

52

53

54 **Abbreviations (footnote on first page)**

55	H/LC	Habitat/Land-Cover
56	RF	Random Forest
57	FRF	Flat Random Forest
58	HRF	Hierarchical Random Forest
59	OoB	Out of Bag
60	H.Step	Hierarchical Stepwise Majority Rule
61	H.Mult	Hierarchical Multiplicative Majority Rule
62	FAO-LCCS	Food and Agriculture Organisation's Land Cover Classification System
63	EUNIS	EUropean Nature Information System habitat classification system
64	Hie.F	Hierarchical F measure

65

66 **1. Introduction**

67 Human-mediated changes in the distribution of habitats and land-cover types are one of the main
68 drivers of the global biodiversity crisis. Consequently, providing reliable Habitat/Land-Cover (H/LC)
69 maps for various conservation related issues is of high priority. For example, H/LC maps are used as
70 input layers for species distribution models (Carlson et al., 2014; Coops et al., 2016; Thuiller et al.,
71 2004) or as obligatory background layers for conservation of umbrella species with well-defined
72 habitat requirements (Li and Pimm, 2016; Murphy and Noon, 1992). Furthermore, H/LC maps are
73 fundamental for mapping ecosystem services (e.g., Koschke et al., 2012) and for natural capital
74 assessments (Brown et al., 2016). Finally, in many cases, the habitats themselves are targeted for
75 conservation and management. For example, as part of the EU Habitat Directive (EU, 2007), all
76 member states of the European Union are required to periodically produce H/LC maps and use the
77 maps for change detection and conservation status assessment. Hence further developing our ability
78 to produce H/LC maps at fine thematic and spatial resolution over wide extents is essential for
79 effective conservation, planning, monitoring, reporting and management of natural resources. As a
80 consequence, there has been a recent surge of methodological and conceptual developments in the
81 field of H/LC classification (Blaschke, 2010; Corbane et al., 2015; Lu and Weng, 2007; Lucas et al.,
82 2015; Lucas et al., 2011; Myint et al., 2011; Tso and Mather, 2009; Xie et al., 2008).

83 In recent years the usage of machine-learning algorithms has become increasingly popular (Belgiu
84 and Drăguț, 2016) as these machine-learning algorithms are efficient at identifying complex
85 classification rule sets, thus potentially providing accurate classification outputs with relatively little

86 investment of time and effort. In many cases, the H/LC classified by machine-learning algorithms rely
87 on a pre-defined national or international classification schemes (e.g., CORINE), to allow a common
88 language of communication between scientists, management agents and policy makers. Most
89 classification schemes adopt a hierarchical, tree-like structure due to several advantages of such
90 structures. Firstly, classes within a hierarchical classification scheme can be grouped into more
91 abstract classes based on semantic similarity criteria, i.e., a hierarchical H/LC class set comprises
92 several semantic granularities. Secondly, a hierarchical H/LC class set can be applied to a variety of
93 spatial scales (each spatial scale requiring the selection of a scale-specific semantic granularity). The
94 former characteristic is particularly useful to meet the minimum required accuracy standard when a
95 specific subclass accuracy is below this standard (Congalton, 1991) and/or when it is difficult to
96 differentiate between subclasses at a given spatial scale. For example, the European Nature
97 Information System habitat classification scheme (EUNIS) has a tree-like structure with up to eight
98 hierarchical levels, containing a total of 5282 habitat classes (at all levels). CORINE LC has three
99 hierarchical levels, with a total of 44 LC classes. Similarly, classification schemes invented ad-hoc for
100 more local studies may also have a hierarchical structure (e.g., Haest et al., 2017). However, most
101 machine learning algorithms follow a flat classification approach (sensu Silla and Freitas, 2011) in
102 which all H/LCs are classified simultaneously in a 'one-against-all' approach. In other words, machine
103 learning algorithms ignore information on the thematic hierarchy that forms the conceptual basis of
104 most classification schemes. Interestingly, many knowledge-based classifiers follow a top-down
105 approach, in which experts first provide rules (e.g., spectral) that separate general H/LC classes from
106 one another, and then move down the thematic tree while providing more specific rules for more
107 specific H/LC categories (e.g., Lucas et al., 2011; Lucas et al., 2007).

108 There are several reasons why incorporating such hierarchical information into the analytical
109 pathway may be beneficial. First, the rule-sets produced by most machine learning algorithms are a
110 'black-box' to the users because of their size and complexity. It is therefore, very difficult to
111 understand or visualise what variables are important in distinguishing between specific sets of
112 habitats. A hierarchical approach may shed some light into the 'black-box' by providing information
113 on variable importance in various locations along the class hierarchy. Second, habitats that are
114 thematically close to one another are not necessarily ecologically/spectrally similar. For example, a
115 forest and grassland may both be listed under the thematic group of 'non-crop' habitats while a
116 wheat field will occur under the thematic group of 'crops' habitats. However, ecologically and
117 spectrally, the grassland may resemble the wheat field more than the forest. A flat classification
118 approach ignores the thematic proximity altogether, while a hierarchical approach will first invest
119 considerable effort in distinguishing 'crop' from 'non-crop', thus potentially preventing confusion

120 between spectrally similar yet thematically distant habitats. Third, if the number of habitats is large,
121 the flat approach may not be able to deal with the complexity of the thematic data, while a
122 hierarchical approach could break the problem into manageable portions by partitioning the
123 feature-space of each group into lower dimensions.

124 Finally, it has been shown that incorporating the hierarchical structure into the modelling
125 framework can increase model accuracy (Thoonen et al., 2013). More specifically, Silla and Freitas
126 (2011) found that various hierarchical approaches tended to increase model accuracy in a wide
127 range of classification problems, especially when misclassifications are weighted by their distance
128 along the classification hierarchy (Kiritchenko et al., 2005). Such hierarchical measures of accuracy
129 acknowledge that not all misclassifications are as critical as the others, e.g., misclassifying one
130 broadleaved-woodland habitat as an alternative closely-related woodland type is arguably a less
131 critical mistake than misclassifying it as a grassland or saltmarsh. In addition, flat classification
132 models only provide performance measures for the entire model or at the H/LC level (i.e., user and
133 producer accuracies). Hierarchical classification models provide the same information with
134 additional accuracy for each local model. That is, the hierarchical approach also provides accuracy
135 for sets of H/LCs that share a common ancestor along the class hierarchy. This information may be
136 crucial for decision makers that may be less interested in the overall accuracy of a map and more by
137 its ability to provide reliable information on sets of H/LCs they care most about (e.g., how well does
138 this model classify non-crop habitats?).

139 We are aware of only a few published manuscripts that focused on hierarchical, machine-learning
140 based classification methods in the remote-sensing literature. Melgani and Bruzzone (2004) found
141 that several support-vector-machine based hierarchical models outperformed flat models when
142 classifying 9 land-use classes in northwest Indiana. Thoonen et al. (2013) found that a tree-structure
143 Markov random field (TS-MRF) method, which captures the hierarchical thematic structure as well as
144 contextual information, outperformed flat classification methods for heathland areas in Belgium.
145 O'Connell et al. (2015) accounted for spatial hierarchy (nested objects) and thematic hierarchy (2
146 levels). They reported slightly better classification outcomes (compared to a flat approach) when the
147 probabilities from a Random Forest (RF; Breiman, 2001) model trained at the top level of the
148 thematic hierarchy were included as predictors of RF models trained at the lower level of the
149 thematic hierarchy. Pena et al. (2014) compared flat and hierarchical approaches (based on 4
150 different algorithms) for mapping cropland areas and found that the flat approach was slightly
151 outperformed by a support-vector-machine based hierarchical model, which fitted a local classifier
152 per parent node. They also found the hierarchical approach increased the minimum sensitivity at the
153 crop level. Finally, Haest et al. (2017) applied an hierarchical classification along four thematic levels

154 when classifying heathland vegetation types for conservation status assessment. They followed a
155 top-down approach such that the class selected for a given pixel in level 2 of the hierarchy could only
156 be one of the children classes of the class selected in level 1 (with similar rules for levels 3 and 4).
157 Haest et al. (2017) observed higher accuracies for the hierarchical approach compared to a flat
158 approach.

159 In this paper we introduce a novel application of hierarchical classification based on the RF
160 algorithm, which accounts for the pre-defined hierarchical structure of classification schemes. The
161 application is available for use in a new R package, entitled '*HieRanFor*'. We tested the hierarchical
162 approach in a NATURA 2000 study site from Italy, using three different classification schemes. Our
163 main aim is to compare the performance of the hierarchical and flat approaches and to explore if the
164 variables identified as important at various locations along the class hierarchy provide meaningful
165 ecological knowledge of the system.

166 **2. Flat and Hierarchical Random Forest**

167 RF is a widely used and well-known classifier in remote sensing (Belgiu and Drăguț, 2016; Bradter et
168 al., 2011). When applied in the flat 'one against all' approach (Flat Random Forest, hereafter, FRF)
169 the model uses a set of training cells and relevant explanatory variables to 'learn' the rules that
170 distinguish one H/LC class from another (Figure 1A,B). The learning procedure is based on fitting a
171 'forest' of classification trees (Figure 1B). Each tree differs from others in the division of the training
172 data to an 'in bag' set (used for tree growing) and an 'Out of Bag' set (OoB, not used for tree
173 growing) and in the usage of variables during the tree growing procedure (Figure 1A). Each tree
174 assigns each OoB case to a single H/LC class. The results are then usually translated into a soft
175 classification output, by estimating for each case, the proportion of OoB votes that assigned the case
176 to any of the H/LC classes (Figure 1B). The reliance on OoB votes and the uniqueness of the
177 constituting trees result in the method being robust to over-fitting. In addition, as prediction are
178 only provided from OoB votes, performance can be assessed without setting aside a considerable
179 portion of the dataset as an external validation set.

180 The Hierarchical Random Forest (HRF) takes similar inputs as FRT along with additional
181 information on the thematic class hierarchy (Figure 1C). Then, RF is used as the local classifier at
182 every internal node of the class hierarchy that has at least two child nodes. The training set for each
183 RF consists of the training cases of all descendants of the parent node. For example, in Figure 1C,
184 training cases representing both H4 and H5 are used in classifier C1 to represent H1. The same cases
185 are also used in local classifier C2 to represent H4 or H5, respectively. However, H4 and H5 training
186 cases are not used in classifier C3 which aims to separate H6, H7 and H8 from one another.

187 After training all the local RFs, it is possible to predict the proportion of votes each case got for
188 each class in each local classifier. Similar to FRF, it is important to ensure the usage of OoB trees
189 when predicting the proportion of votes, otherwise we may overestimate the performance of the
190 model and may be prone to overfitting. For example, a case from the H4 class of figure 1 is included
191 in the training data of classifiers C1 and C2. For these two classifiers, we would only use the votes of
192 trees in which the case was not randomly assigned to the in-bag set (tree training set, Fig. 1A). In
193 classifier C3, which only included cases from H6, H7 and H8 in the training data, all trees can be used
194 for predicting the H4 case probabilities. If a case is never predicted by a classification tree in which it
195 is within the 'in-bag' subset, HRF should retain the robustness of FRF to over-fitting. Similarly, as
196 accuracy is always based on OoB votes, there is no need for an external validation set. Nonetheless,
197 similar to FRF, it is possible to run additional cases through the HRF model for prediction and
198 producing maps or for additional assessment of performance using an external validation set. If the
199 hierarchical model is used to predict for new data, all trees from all local classifiers can be used.

200 Both FRF and HRF produce for each case the proportion of OoB votes for each H/LC. To assess
201 the performance (e.g., kappa) of the model and to produce maps, a single habitat needs to be
202 selected according to the proportion of OoB votes (i.e., translating the 'soft' probabilities to a 'crisp'
203 classification). For FRF, the flat majority rule is usually applied by selecting for a focal case the H/LC
204 that received the highest proportion of OoB votes. For HRF, there are two different majority rule
205 options – stepwise majority (H.Step) and multiplicative majority (H.Mult, Figure 1C). In stepwise
206 majority rule, the flat majority rule is applied to each local classifier, and starting from the tree root,
207 the selected H/LC is followed until a terminal H/LC (with no descendants lower down the hierarchy)
208 is reached. In the case presented in Figure 1C (blue), H3 received the highest proportion of OoB
209 votes at classifier C1 (0.5 versus 0.1 and 0.4), while H6 received the highest proportion of OoB votes
210 in classifier C3 (0.6 versus 0.3 and 0.1). In the multiplicative majority rule, the proportion of OoB
211 votes are multiplied along every path from the tree root until it reaches the terminal H/LC, and the
212 flat majority rule is applied on the multiplicative proportions. In figure 1C, the multiplicative votes
213 identify H4 (red) as the most probable H/LC. Interestingly, the multiplicative proportion of OoB votes
214 are comparable to the proportion of OoB votes generated by FRF in a sense that both sum to 1 over
215 all terminal nodes.

216 **3. Methods**

217 *3.1 Study Site*

218 We have focused our analysis on Le-Cesine, Italy (IT1 – SCI IT9150032; SPA IT9150014) – a Natura
219 2000 site and one of the oldest protected areas in Puglia. Le-Cesine covers an area of about 2148 ha.

220 The site is characterized by a high diversity in land-cover, habitat and vegetation types. The coastal
221 wetland is characterized by a system of lagoons, ponds and marshes where several types of
222 helophytic vegetation (reeds, sedges and rushes communities) are widespread. Further inland, the
223 woody vegetation is composed by a mosaic of *Pinus halepensis* stands and different types of
224 Mediterranean maquis and garrigues, while the agricultural areas are mainly composed of olive
225 groves. The site is affected by marine erosion of the sandbank, resulting in reduction and
226 fragmentation of the typical dune habitat types, as well as salinization of the lagoons and the related
227 environments.

228 3.2 *Ground-Truth Data and Classification Schemes*

229 The entire extent of Le-Cesine was ground-truthed in earlier projects to two LC classifications
230 schemes and one H classification scheme. The two LC classification schemes included CORINE
231 (Bossard et al., 2000) and the UN Food and Agriculture Organisation's Land Cover Classification
232 System (FAO-LCCS, Di Gregorio and Jansen, 2005), while for the H classification scheme we relied on
233 EUNIS (Davies and Moss, 2002). The three classification schemes differ from one another in the
234 number of internal and terminal nodes, in the number of hierarchical levels, and in the number of
235 local classifiers required to run the HRF analysis.

236 A pre-existing validated LC map (scale 1:5000) in CORINE Land Cover was available from a
237 previous Interreg (Nat Info) project. The selection of an appropriate LC classification system for
238 habitat mapping applications is a crucial issue. The FAO-LCCS (Di Gregorio and Jansen, 2005) has
239 been considered as an appropriate and user-friendly framework for long-term monitoring of the
240 conservation status of habitats. LCCS allows the finest discrimination of natural and semi-natural
241 types with respect to other widely used LC taxonomies (Tomaselli et al., 2013). Thus, CORINE classes
242 were first converted to FAO-LCCS classes according to the LCCS2 software (Di Gregorio and Jansen,
243 2005) and the semantic heterogeneity issues were addressed through the expert knowledge of
244 people with long experience on the monitoring of the study site. Then, FAO-LCCS classes were
245 translated to EUNIS habitat classes by integrating the environmental attributes, i.e. lithology, soil
246 group, soil-surface aspect and water quality which are the auxiliary data that can be used for habitat
247 discrimination.

248 This information was provided from a previous INTERREG (Nat Info) project. As well known
249 (Adamo et al., 2016; Tomaselli et al., 2013) translating LC to habitats classes may include one-to-
250 many relations, so prior-knowledge from botanists was needed for selecting the environmental data
251 useful to discriminate the different habitat classes that may correspond to a specific LC class. These
252 environmental attributes help to resolve, in most cases, the challenge of one-to-many relationships

253 between LC and habitat classes. Field surveys were then carried out in 2011-2013 to validate both
254 FAO-LCCS and habitat classes, according to the EUNIS report (Davies and Moss, 2002) obtained from
255 each CORINE LC class and to select reference samples for training/testing the RF outputs. As a result,
256 the number of input CORINE classes is lower than the number of FAO-LCCS classes and EUNIS classes
257 (see Figure 2 in Tomaselli et al., 2013). A random sampling design within a 250 m cell regular grid, in
258 turn nested within a 1 km cell standard regular grid (INSPIRE), was selected; within this grid, 50
259 circular 50 m radius vegetation plots, randomly distributed throughout the site and covering all the
260 habitat types (according to EUNIS), were recorded and mapped. For each of these points,
261 information was collected on vegetation composition and structure, crop cover and habitat type
262 (Tomaselli et al., 2016). Such information, geocoded by GPS, was integrated into a GIS geo-database
263 using ArcGIS 9.2.

264 The CORINE classification scheme (Figure 2) for our study system contains 31 land cover classes
265 organized in three hierarchical levels, with 14 of the classes being terminal (i.e., not further
266 described as any other more detailed classes). A total of 8 local RF classifiers are required to fit the
267 HRF model according to the CORINE classification scheme. The FAO-LCCS classification scheme
268 (Figure 3) contains 49 land cover classes (18 terminal) organized along 6 hierarchical levels and
269 requiring 8 local RF classifiers. The EUNIS classification scheme (Figure 4) contains 60 habitats (23 of
270 which are terminal) organized along 5 hierarchical levels and requiring 11 local RF classifiers. Thus,
271 we can explore here the effect of the shape (number of hierarchical levels and distribution of nodes
272 along the tree structure) of the hierarchical classification scheme on the performance of the
273 different classification methods, in the same system using the exact same pixels and explanatory
274 variables.

275 3.3 Data Preparation and Explanatory Variables

276 Detailed information on the processing procedure of the images, on the derived indices and on each
277 of the other explanatory variables is found in supporting information S1. We based our analysis on
278 61,453 cells, at 10x10 m resolution covering the entire extent of the study-site. We used a 10x10 m
279 resolution, as preliminary analysis revealed it to produce similar results as finer resolutions (2x2 m
280 and 5x5 m) with considerably lower running time, while still providing detailed enough information
281 for management in this study system. For each cell we calculated 35 explanatory variables. The first
282 12 variables were radiometrically calibrated reflectance values of two Very High Resolution (VHR)
283 remotely sensed images: 1) a multispectral WorldView-2 image with 8 spectral bands at 2m spatial
284 resolution taken on October 9th 2010; 2) a multispectral Quickbird image with 4 spectral bands at 2.4
285 m spatial resolution taken on June 4th 2009. We applied radiometric correction in accordance with

286 O'Connell et al. (2013) where top-of-atmosphere reflectance followed by dark object subtraction
287 correction was completed. The images were geometrically synced using 221 ground control using a
288 direct linear transform model, giving a root-mean-squared-error of 0.6. Resampling was then done
289 on both images to grids of 10 m using a cubic convolution interpolation, ensuring that both images
290 were radiometrically and spatially aligned.

291 Five vegetation indices were also derived from each image; Difference Vegetation Index (DVI),
292 Normalised Difference Vegetation Index (NDVI), Principal Components Analysis (PCA), Soil Adjusted
293 Vegetation Index (SAVI) and Atmos Resistant Vegetation Index (ARVI). For the WorldView-2 image
294 we added two additional derived indices- the WorldView-2 Soil Index (WVSI) and WorldView-2
295 improved Vegetation Index (WVVI) based on the additional red edge, coastal, yellow and near-Infra
296 Red2 bands. For each of the four bands of the Quickbird image we further calculate the Tasselled
297 Cap Index. Therefore a total of 16 spectrally derived variables were created from the two satellite
298 datasets. Furthermore, we estimated 7 environmental variables including elevation, based on a 10 m
299 resolution DEM model (Tarquini et al., 2007; Tarquini et al., 2012) and plant height layers obtained
300 from a LiDAR campaign (canopy height model). The last five variables were categorical
301 environmental variables covering the lithology (3 classes), soil group (6 levels), Soil-Surface aspect
302 (Soil Bare, 6 levels), water quality (6 levels) and cadastral information (16 levels). Such information
303 were collected from local authorities' archives and public repositories.

304

305 3.4 *Creating Datasets for Classifications*

306 Both the FRF and HRF require a training set containing cases (pixels or object) from each of the H/LC
307 classes. To create a dataset, we first randomly selected from each H/LC 100 cases for the training
308 set. The remaining cases were used as an external, independent validation set. If the total number of
309 available cases for a given H/LC was lower than 100, we included all its cases in the training data. We
310 repeated this procedure 15 times per classification scheme (45 datasets all together), with each
311 dataset containing both a small training subset and a much larger independent validation subset.

312 3.5 *Fitting the FRF and HRF Models for a Single Dataset*

313 As part of the FP7 European project EU BON (Hoffmann et al., 2014) we created the '*HieRanFor*' R
314 package which provides full functionality of applying the above framework (development version is
315 available at <https://r-forge.r-project.org/projects/hie-ran-forest/>, while a working version is
316 available as SI for this publication). The '*HieRanFor*' uses the RF algorithm as implemented in the
317 '*randomForest*' package (Liaw and Wiener, 2002) in R (R Core Team, 2016). With the '*HieRanFor*'

318 package, users can fit an HRF, predict the proportion of OoB votes in each local classifier (for both
319 the training data and for external validation data) and translate the proportion of OoB votes to a
320 crisp class using either stepwise or multiplicative majority rule. The *'HieRanFor'* package further
321 includes functions to estimate the performance of the model using either flat or hierarchal indices.
322 Finally, the package allows extraction of variable importance values for each class in each local
323 classifier.

324 We first used the training data to fit a FRF model and an HRF model. For the FRF model we used
325 exclusively functions available in the *'randomForest'* package for modelling, extracting the
326 proportion of OoB votes for the training set, and predicting the proportion of votes for the validation
327 set. For the HRF model we first ran the model using the *'RunHRF'* function of *'HieRanFor'*. Then, we
328 used the function *'PerformanceHRF'* to extract the proportion of OoB votes for each local classifier,
329 as well as the multiplicative proportion of OoB votes, and the *'PerformanceNewHRF'* function for
330 predicting the proportion of votes for the validation dataset. At this stage, we had for each case of
331 the training and validation set the proportion of OoB votes required to translate the soft
332 classification into a crisp one.

333 In FRF, we translated the soft classification into crisp by employing a simple majority rule, i.e.,
334 selecting for each case the H/LC that received the highest proportion of OoB votes. For HRF, we used
335 both the stepwise majority rule and the multiplicative majority rule (see above and in Figure 1C). The
336 crisp H/LC for each case were then used to assess the performance of the models separately for the
337 FRF, the stepwise HRF, and the multiplicative HRF. We estimated performance separately for the
338 training and validation set, using the overall accuracy (based on the diagonal of the error matrix),
339 Cohen's Kappa (Cohen, 1968), and the Hierarchical F measure (Kiritchenko et al., 2005), for a total of
340 18 performance values for each dataset. The unweighted kappa and accuracy were computed using
341 the *confusionMatrix* function of the *'caret'* R package (Kuhn et al., 2016). For the Hierarchical F
342 measure (Hie.F), we implemented within *'HieRanFor'* a function that calculates the Kiritchenko et al.
343 (2005) index directly from the confusion matrix. Additional information on Hie.F can be found in
344 supporting information 2.

345 3.6 Variable Importance

346 In FRF, variable importance is quantified as the mean decrease in accuracy of the entire model when
347 the values of a focal variable are permuted, and all other values remain unchanged. In HRF,
348 variable importance is estimated separately for each variable in each local classifier, thereby
349 allowing more information on why certain variables are more important than others. We explored

350 the variable importance separately for each classification scheme by taking the mean over all 15
351 runs.

352 3.7 Statistical Analysis

353 After fitting the FRF and HRF for the 45 datasets following the procedure outlined above, we
354 explored the models performance using mixed-effect models. We used one of the model
355 performance indices (kappa, accuracy or Hie.F measure) as the dependent variables. For fixed
356 effects, we used the data type (2 levels- training or validation), the crisp rule (3 levels: simple
357 majority for the FRF, stepwise for the HRF and multiplicative for the HRF) and their interaction. For
358 random effects we used the data type, grouped according to the model run (i.e., the specific set of
359 training and validation datasets created randomly -- a paired design). We used this model since it
360 received the lowest AICc (Akaike Information Criteria corrected for small sample size values) from a
361 set of 9 nested mixed effect models (supporting information S3). We repeated the analysis
362 separately for each of the three H/LC classification schemes.

363 4. Results

364 In each run, we have used a total of 1333, 1661 and 1980 pixels in the training set for CORINE,
365 FAO-LCCS and EUNIS, respectively -- only 2.2-3.2% of the study area's 61453 pixels (supporting
366 information S4 for breakdown of the training set to classes and for examples of several confusion
367 matrices). Despite the small training set, all methods in all classification schemes produced maps
368 that are very similar to the observed map, both for the small training set and the large validation set
369 (see examples in Figures 2-4). In all 3 classification schemes, we found the environmental variables
370 to have considerably higher variable importance values relative to the spectral reflectance values
371 and the vegetation indices derived from the remotely-sensed images (Figure 2-4, lower left panels).
372 Among the environmental variables, water quality, soil group, soil bare and lithology showed the
373 highest variable importance values, followed by the cadastral and elevation variables, and finally
374 canopy height. This pattern was observed for both the FRF and HRF approaches in all 3 classification
375 schemes. However, the HRF approach provided more detailed information on the importance of
376 each variable in each local classifier, flagging different variables as important in different locations
377 along the hierarchy, even if they had relatively low importance when averaged across all categories
378 using FRF. For example, elevation was identified as the most important variable in local classifier C.8
379 for CORINE (Figure 2). Similarly, the cadastral variable received relatively low relative importance
380 score in the flat model, but it was identified as the most important variable in one of the local
381 classifiers of the FAO-LCCS (C.4 in Figure 3). This may be due to the fact that cadastral allows the
382 discrimination of barren land from artificial structures for objects characterized by similar spectral

383 signature, when additional context-sensitive features are not considered (e.g., texture). For EUNIS,
384 water quality was identified as having very low importance for some classifiers but was very
385 important in several others, including being the only important variable for separating the two
386 children nodes of habitat A2.52 (Figure 4, C.10). This is justified by the importance of water salinity
387 for discriminating EUNIS habitats such as C3.421 (fresh water) and A2.51 (salt water) corresponding
388 to the same FAO-LCCS class (A24/A2.A5.E7, i.e., inland water habitats (see Table 6 in Tomaselli et al.,
389 2013).

390 The overall performance of both the FRF and HRF was very high. Kappa values were above 0.9 for
391 CORINE and FAO-LCCS, with slightly lower values for EUNIS (Figure 5, middle column). Accuracy
392 values were above 0.9 for CORINE and FAO-LCCS, i.e., more than 90% of the cells were correctly
393 classified (Figure 5, left column). Accuracy for EUNIS was slightly lower, but still above 0.86 in all
394 runs. Similar high values were also observed for the Hie.F measure (Figure 5, right column). The
395 variance between the 15 runs for a given combination of performance index and classification
396 scheme was very low, suggesting that the choice of training set had little effect on the results. Given
397 the overall high performance values, it was hard to detect performance differences between the FRF
398 and HRF approaches (Figure 5). In CORINE we observed a slight decrease in all three performance
399 measures in the two HRF approaches relative to the FRF, for both the training and validation
400 datasets. For EUNIS, FRF performed slightly better in the training set, yet the two HRF approaches
401 slightly outperformed the FRF in the validation set. A similar pattern regarding the hierarchical
402 stepwise approach was also observed for the FAO-LCCS classification scheme. When comparing the
403 medians, while the FRF outperformed the two HRF in all 9 cases in the training set (3 classification
404 schemes \times 3 performance indices), the HRF outperformed the FRF in 6 out of 9 cases in the
405 validation set. The hierarchical multiplicative approach outperformed the hierarchical stepwise
406 approach in EUNIS and CORINE, while the hierarchical stepwise approach outperformed the
407 hierarchical multiplicative approach in FAO-LCCS.

408 In the validation datasets, the pairs of H/LC that were responsible for the highest number of
409 classification errors in each of the three schemes are summarized in table 1. In general, the main
410 source of confusion over all schemes is amongst the olive groves and conifer plantations. The
411 inclusion of more images covering the seasonal cycle and/or the introduction of context-sensitive
412 features related to agricultural practices may improve the discrimination (e.g., olive trees are
413 organized in parallel rows at regular distances). These habitats were either misclassified as one
414 another, or as road or fields. In CORINE, confusion between *Coniferous Forest* (3.1.2) or *Olive Groves*
415 (2.2.3) on one side and *Permanently Irrigated Lands* (2.1.2) or *Road and Rail Networks and*
416 *Associated Land* (1.2.2) on the other side accounted for 71.5%, 76.9% and 75.0% of all classification

417 errors for the flat, H.Mult and H.Step approaches, respectively. In FAO-LCCS, confusion between
418 *Plantations: needle-leaved evergreen tree crops- monoculture + rainfed*
419 (A11_A1.B1.A8.A9.B3.W7.C1.D1) or *Orchards: broad-leaved evergreen tree crops- monoculture +*
420 *rainfed* (A11_A1.B1.A7.A9.B4.W8.C1.D1) on one side and *Fields of irrigated no graminoid crops + one*
421 *additional crop* (A11_A3.A5.B2.C2.D3) or *Paved roads* (B15_A1.A3.A7.A8) on the other side
422 accounted for 55.8%, 55.5% and 47.2% of all classification errors for the flat, H.Mult and H.Step
423 approaches, respectively. In EUNIS, confusion between *Native Conifer Plantations* (G3.F1) or *Olea*
424 *europaea Groves* (G2.91) on one side and *Arable Land with Unmixed Crops Grown by Low-intensity*
425 *Agricultural Methods* (I1.3) or *Road Networks* (J4.2) on the other side accounted for 40.3%, 40.0%
426 and 38.4% of all classification errors for the flat, H.Mult and H.Step approaches, respectively. In
427 addition, for EUNIS (and to a lesser extent, for FAO-LCCS) misclassification of different types of salt
428 marshes as one another also contributed considerably to the overall error rate. These errors were
429 avoided in CORINE due to the lower thematic resolution.

430 **5. Discussion**

431 *5.1 Performance of the Flat and Hierarchical Models*

432 In this manuscript we introduce the HRF: a machine learning classification algorithm that accounts
433 for the pre-defined hierarchical structure of user-defined classification schemes. Our main objective
434 was to explore whether this novel hierarchical approach provides better performance relative to the
435 commonly used flat classification approach. We further explored whether the additional information
436 provided by the HRF approach could shed light into the 'black-box' of RF. A single study site was
437 used, for which we had complete independent coverage data based on two different LC and one H
438 classification schemes. In general, we have found that for all three classification schemes, both the
439 FRF and the HRF were able to predict the observed H/LC with high accuracy (Figure 5), despite using
440 as low as 2.2% of the study site as training cells. This was observed not only in the training sets, but
441 also in the external, much larger and independent validation sets.

442 Given the very high performance of both FRF and HRF, it is difficult to identify explicit
443 differences between the two approaches. The difference in performance between the FRF, H.Mult
444 and H.Step were very small. However there were some differences, with the flat approach best for
445 CORINE, the H.Step best for FAO-LCCS and the H.Mult best for EUNIS (Figure 5). Although a single
446 case-study is not enough for generalization, these differences may be attributed to differences in
447 class hierarchy. The hierarchical approaches outperformed the flat one in the two classification
448 schemes that have a complex hierarchical structure, with more nodes and more levels. In addition,
449 the H.Step outperformed the H.Mult in FAO-LCCS, perhaps since FAO-LCCS classifies only two classes

450 in the local classifier closest to the tree root (C.1 in Figure 3) while EUNIS already has ten classes in
451 C.1 (Figure 4). The small number of classes in the first classifier provides an advantage to the
452 stepwise approach, with the higher weights it gives to the early classifier affecting the overall
453 performance. Thus, based on this single case-study we suggest use of the flat approach when the
454 hierarchy is relatively simple, the H.Step approach when the hierarchy is more complex but with
455 simple early classifiers and the H.Mult when the hierarchy is complex right from the start.

456 In fact, as clearly noted by Haest et al. (2017), the main advantage of the H.Step approach may
457 be its ability to separate thematically distinct but spectrally similar classes, assuming the classifier
458 close to the tree-root performs well. However, when this classifier does not perform well, the errors
459 will be carried down the class structure with no option to undo them. The H.Mult approach does not
460 give higher weights to classifiers close to the tree root and thus errors made in certain classifiers
461 have only a limited effect on the overall predicted probabilities. Instead, the probabilities are
462 multiplied along the tree root along every path. Thus, probabilities of classes that have more
463 children nodes further down the hierarchy may be diluted relative to other classes that terminate
464 closer to the tree root. This can potentially bias the method toward shallower classes. If such a bias
465 affected our results, we would expect the H.Mult to classify more cases to classes that are closer to
466 the tree root than the flat and H.Step approaches. We have not found any evidence for such a bias
467 when comparing the number of cases classified to each level in the two classification schemes that
468 have terminal classes at different levels (FAO-LCCS and EUNIS).

469 5.2 *Types of Misclassifications*

470 Both FRF and HRF made most misclassification errors when trying to classify two main H/LC types
471 – the conifer plantations and the olive groves (table 1). In these H/LC types, the actual cover of trees
472 is considerably lower than 100%. The gaps between the trees are covered with bare soil or with
473 green/dry annual vegetation (depending on the intensity of grazing, the agricultural practice and the
474 time of the year). Thus, the main source of confusion was treating the gaps' bare soil area as roads
475 or the gaps' annuals areas as agricultural fields. These misclassifications most likely represent cases
476 that are thematically distant yet spectrally similar. Our general expectation was that HRF will make
477 fewer such errors than FRF, since local classifiers closer to the tree root would already separate the
478 H/LC at early stage. We further expected the stepwise approach to be better than the multiplicative
479 approach in that respect, since the stepwise approach gives higher weights to local classifiers closer
480 to the tree root while the multiplicative approach gives equal weights to all local classifiers. Indeed,
481 we observed such a pattern for EUNIS, with 2669 vs. 2830 misclassifications of the above types for
482 H.Step and FRF, respectively. Similarly, for the FAO-LCCS there were 2291 vs. 2851 such

483 misclassifications for the H.Step and FRF, respectively. For the FAO-LCCS, the difficulty in classifying
484 gaps resulted for H.Step mainly in confusion between the orchards and plantations, that are
485 classified at the same local classifier (C.4 in Figure 3). Interestingly, in CORINE, the flat approach
486 outperformed the H.Step (2793 vs. 3178 misclassifications, respectively).

487 5.3 Variable importance

488 In all three classification schemes, the environmental variables were more important than the
489 raw spectral variables or the variables derived from the remotely sensed images (lower left panels in
490 figures 2-4). Perhaps, the poor performance of the spectral variables is due to a scale issue; spectral
491 variables in high resolution images change over very short distances on the ground (even after
492 resampling) but the environmental variables were based on relatively coarse resolution data and
493 may have tied in better with the spatial scale of the classification. This can be a common issue in
494 complex ecosystems or habitats. In addition, the remote sensing data were only from two dates (one
495 of which was approaching winter leaf-off) and the spectral variables may increase in importance if
496 more dates were available at key phenological stages in the season.

497 Nonetheless, the HRF approach also provided more detailed information on the importance of
498 each variable in each local classifier. For example, in the EUNIS classification scheme (Figure 4) the
499 local classifier C.10 reveals that distinguishing between the *Mediterranean Juncus maritimus* and
500 *Juncus acutus saltmarshes* (A2.522) and *Mediterranean saltmarsh scrubs* (A2.526) requires
501 information on water quality (salinity). However, distinguishing between *Marine saline beds of*
502 *Phragmites australis* (A2.53C) and *Geolittoral wetlands and meadows: saline and brackish reed, rush*
503 *and sedge stands* (A2.53D) in local classifier C.11 requires information on both water quality and soil
504 groups, perhaps because A2.53C has a much wider ecological range with its distribution ranging
505 from the dune area to the inland area. In FAO-LCCS (Figure 3) the first local classifier (C.1) reveals
506 that none of the available explanatory variables are extremely important in distinguishing between
507 primarily vegetated areas (habitat A) and primarily non-vegetated areas habitat (B). In the FAO-LCCS
508 (Figure 3), the explanatory variables of cadastral information and the canopy height model were not
509 identified as important by the flat model. However, these are the main explanatory variables that
510 distinguish between the three descendent H/LC of *Cultivated and Managed Terrestrial Areas* (A11) in
511 local classifier C.4. Similarly, in the CORINE classification scheme (Figure 2), the elevation variable is
512 not identified as important by the flat model, yet the HRF analysis reveals it to have a relatively high
513 effect on accuracy in at least 2 local classifiers (e.g., C.2 and C.8).

514 5.4 The Potential of Hierarchical Classification Models

515 The hierarchical approach we tested here is only one of several potential hierarchical approaches
516 that have been used for classification problems in other research areas. In fact, in their review, Silla
517 and Freitas (2011) list several types of hierarchical models including:

- 518 • *Local classifier per parent node* – a local classification model is fitted in every internal node that
519 has more than one sibling. This is the approach followed here and in Pena et al. (2014).
- 520 • *Local classifier per node*- a binary model is fitted to each internal and terminal node, e.g., 9 local
521 models in the example of Figure 1, each for one habitat (H1-H8).
- 522 • *Local classifier per level* – A local model is fitted for each level of the hierarchy, classifying all the
523 habitats in the level. In the example of Figure 1, two classification models would be fitted, one
524 for classifying H1, H2 and H3 and one classifying H4, H5, H6, H7 and H8.
- 525 • *Global classifier* – a single model is fitted which simultaneously captures the entire class
526 hierarchy. As far as we know, this approach has never been tested using machine learning
527 algorithm for remote-sensing. Interestingly, knowledge-based classifiers, which tend to follow a
528 hierarchical classification structure, are the closest example we could think of.

529 It is probable that different approaches may be more suitable for particular applications. For
530 example, the local classifier per node approach may be most suitable for situations in which not all
531 cases in the training data are classified to the lowest level (e.g., some cases in Figure 1C are labelled
532 as H1, without further labelling as either A4 or A5). Alternatively, local classifier per level may be
533 most suitable for cases in which all terminal classes are at the same level (e.g., CORINE), especially if
534 the probabilities of models closer to the tree root are entered as explanatory variables in models
535 lower down the hierarchy (O'Connell et al., 2015). The local classifier per parent node approach may
536 better suit cases in which the hierarchy is more complex (e.g., EUNIS).

537 We can further envision usage of hierarchical models even in the absence of a pre-defined
538 thematic scheme as a way to focus models on sets of classes that are not easily distinguished from
539 one another. In the local classifier per parent approach this can be done by first fitting a flat model
540 to produce a confusion matrix. Then hierarchical cluster analysis may be used to create the thematic
541 hierarchy, such that classes that are commonly confused (e.g. olive and pine plantations) will be
542 closer to each other in the thematic tree. A hierarchical model based on such a tree could then fit
543 specific local classifiers to the most hard-to-distinguish sets.

544 Even if the marginal increase in performance of hierarchical vs. flat models would be low, the
545 more detailed variable importance information of hierarchical models may provide information on
546 the processes that govern the distribution of various habitats. Similarly, the accuracy of specific local
547 classifiers may help identify sets of habitats that none of our predictors can distinguish between.

548 Finally, hierarchical models can be used as an automated tool to lower the costs of knowledge-based
549 classification models, by suggesting variables that are most suitable at distinguishing groups of
550 habitats from one another.

551 **6. Conclusion**

552 This is one of the first applications of a machine learning classification algorithm that accounts for a
553 pre-defined hierarchical thematic structure (but see: Haest et al., 2017; O'Connell et al., 2015; Pena
554 et al., 2014; Thoonen et al., 2013). We did not observe considerable differences in performance
555 between this new hierarchical approach and a typical flat approach. This is in accordance with Pena
556 et al. (2014) that followed a similar local classifier per parent node approach. However, our results
557 do suggest that the best strategy may depend on the complexity of the hierarchical structure.
558 Furthermore, we found that the hierarchical approach provides valuable information on variable
559 importance that shed some light into the usually 'black-box' of machine-learning algorithms. We
560 hope that this work will trigger additional works on the potential of using hierarchical models for
561 H/LC classification.

562

563 **Acknowledgment**

564 We wish to thank the associate editor prof. Jay Gao. We further thank Birgen Haest and an
565 anonymous reviewer for their valuable and constructive comments that improved the quality and
566 clarity of this manuscript.

567 **Funding**

568 This manuscript is a collaboration between two European projects. YG, JO, CJM and WEK were
569 financed by the EU BON project (www.eubon.eu) that is a 7th Framework Programme funded by the
570 European Union under Contract No. 308454. This work was also supported by the Horizon2020
571 ECOPOTENTIAL project: improving future ecosystem benefits through earth observations, grant
572 agreement 641762 (www.ecopotential-project.eu).

573

574 **References**

575 Adamo, M., Tarantino, C., Tomaselli, V., Veronico, G., Nagendra, H., Blonda, P., 2016. Habitat
576 mapping of coastal wetlands using expert knowledge and Earth observation data. *J. Appl. Ecol.*
577 53, 1521-1532.

578 Belgiu, M., Drăguț, L., 2016. Random forest in remote sensing: A review of applications and future
579 directions. *ISPRS-J. Photogramm. Remote Sens.* 114, 24-31.

580 Blaschke, T., 2010. Object based image analysis for remote sensing. *ISPRS-J. Photogramm. Remote*
581 *Sens.* 65, 2-16.

582 Bossard, M., Feranec, J., Otahel, J., 2000. CORINE land cover technical guide-addendum 2000
583 Technical Report, No 40, European Environmental Agency.

584 Bradter, U., Thom, T.J., Altringham, J.D., Kunin, W.E., Benton, T.G., 2011. Prediction of National
585 Vegetation Classification communities in the British uplands using environmental data at multiple
586 spatial scales, aerial images and the classifier random forest. *J. Appl. Ecol.* 48, 1057-1065.

587 Breiman, L., 2001. Random forests. *Mach. Learn.* 45, 5-32.

588 Brown, C., King, S., Ling, M., Bowles-Newark, N., Ingwall-King, L., Wilson, L., Pietilä, K., Regan, E.C.,
589 Vause, J., 2016. Natural Capital Assessments at the National and Sub-national Level. UNEP-
590 WCMC, Cambridge, UK.

591 Carlson, B.Z., Georges, D., Rabatel, A., Randin, C.F., Renaud, J., Delestrade, A., Zimmermann, N.E.,
592 Choler, P., Thuiller, W., 2014. Accounting for tree line shift, glacier retreat and primary succession
593 in mountain plant distribution models. *Diversity and Distributions* 20, 1379-1391.

594 Cohen, J., 1968. Weighted kappa: nominal scale agreement with provision for scaled disagreement
595 or partial credit. *Psychological bulletin* 70, 213-220.

596 Congalton, R.G., 1991. A review of assessing the accuracy of classifications of remotely sensed data.
597 *Remote Sensing of Environment* 37, 35-46.

598 Coops, N.C., Waring, R.H., Plowright, A., Lee, J., Dilts, T.E., 2016. Using Remotely-Sensed Land Cover
599 and Distribution Modeling to Estimate Tree Species Migration in the Pacific Northwest Region of
600 North America. *Remote Sens.* 8.

601 Corbane, C., Lang, S., Pipkins, K., Alleaume, S., Deshayes, M., García Millán, V.E., Strasser, T., Vanden
602 Borre, J., Toon, S., Michael, F., 2015. Remote sensing for mapping natural habitats and their
603 conservation status – New opportunities and challenges. *Int. J. Appl. Earth Obs. Geoinf.* 37, 7-16.

604 Davies, C.E., Moss, D., 2002. EUNIS habitat classification. Final report to the European topic centre of
605 nature protection and biodiversity, European Environment Agency, Swindon.

606 Di Gregorio, A., Jansen, L.J.M., 2005. Land Cover Classification System (LCCS): classification concepts
607 and user manual. Food and Agriculture Organization of the United Nations, Rome.

608 EU, 2007. Habitats Directive. In: Commission, E. (Ed.), Article 10.

609 Haest, B., Vanden Borre, J., Spanhove, T., Thoonen, G., Delalieux, S., Kooistra, L., Múcher, C.,
610 Paelinckx, D., Scheunders, P., Kempeneers, P., 2017. Habitat Mapping and Quality Assessment of
611 NATURA 2000 Heathland Using Airborne Imaging Spectroscopy. *Remote Sens.* 9, 266.

612 Hoffmann, A., Penner, J., Vohland, K., Cramer, W., Doubleday, R., Henle, K., Köljalg, U., Kühn, I.,
613 Kunin, W.E., Negro, J., Penev, L., Rodríguez, C., Saarenmaa, H., Schmeller, D., Stoev, P.,
614 Sutherland, W., Tuama, É.Ó., Wetzell, F., Häuser, C., 2014. The need for an integrated biodiversity
615 policy support process – Building the European contribution to a global Biodiversity Observation
616 Network (EU BON). *Nature Conservation* 6, 49-65.

617 Kiritchenko, S., Matwin, S., Famili, F., 2005. Functional annotation of genes using hierarchical text
618 categorization, In: *Proc. of the ACL Workshop on Linking Biological Literature, Ontologies and*
619 *Databases: Mining Biological Semantics.*

620 Koschke, L., Fürst, C., Frank, S., Makeschin, F., 2012. A multi-criteria approach for an integrated land-
621 cover-based assessment of ecosystem services provision to support landscape planning. *Ecol.*
622 *Indic.* 21, 54-66.

623 Kuhn, M., Contributions from: Wing, J., Weston, S., Williams, A., Keefer, C., Engelhardt, A., Cooper,
624 T., Mayer, Z., Kenkel, B., R Core Team, Benesty, M., Lescarbeau, R., Ziem, A., Scrucca, L., Tang, Y.,
625 Candan, C., Hunt, T., 2016. caret: Classification and Regression Training. R package. version 6.0-
626 73. <https://CRAN.R-project.org/package=caret>.

627 Li, B.V., Pimm, S.L., 2016. China's endemic vertebrates sheltering under the protective umbrella of
628 the giant panda. *Conserv. Biol.* 30, 329-339.

629 Liaw, A., Wiener, M., 2002. Classification and Regression by randomForest. *R News* 2, 18-22.

630 Lu, D., Weng, Q., 2007. A survey of image classification methods and techniques for improving
631 classification performance. *Int. J. Remote Sens.* 28, 823-870.

632 Lucas, R., Blonda, P., Bunting, P., Jones, G., Inglada, J., Arias, M., Kosmidou, V., Petrou, Z.I., Manakos,
633 I., Adamo, M., Charnock, R., Tarantino, C., Múcher, C.A., Jongman, R.H.G., Kramer, H., Arvor, D.,
634 Honrado, J.P., Mairota, P., 2015. The Earth Observation Data for Habitat Monitoring (EODHaM)
635 system. *Int. J. Appl. Earth Obs. Geoinf.* 37, 17-28.

636 Lucas, R., Medcalf, K., Brown, A., Bunting, P., Breyer, J., Clewley, D., Keyworth, S., Blackmore, P.,
637 2011. Updating the Phase 1 habitat map of Wales, UK, using satellite sensor data. *ISPRS-J.*
638 *Photogramm. Remote Sens.* 66, 81-102.

639 Lucas, R., Rowlands, A., Brown, A., Keyworth, S., Bunting, P., 2007. Rule-based classification of multi-
640 temporal satellite imagery for habitat and agricultural land cover mapping. *ISPRS-J. Photogramm.*
641 *Remote Sens.* 62, 165-185.

642 Melgani, F., Bruzzone, L., 2004. Classification of hyperspectral remote sensing images with support
643 vector machines. *IEEE Trans. Geosci. Remote Sensing* 42, 1778-1790.

644 Murphy, D.D., Noon, B.R., 1992. Integrating scientific methods with habitat conservation planning -
645 reserve design for northern spotted owls. *Ecol. Appl.* 2, 3-17.

646 Myint, S.W., Gober, P., Brazel, A., Grossman-Clarke, S., Weng, Q.H., 2011. Per-pixel vs. object-based
647 classification of urban land cover extraction using high spatial resolution imagery. *Remote*
648 *Sensing of Environment* 115, 1145-1161.

649 O'Connell, J., Bradter, U., Benton, T.G., 2015. Wide-area mapping of small-scale features in
650 agricultural landscapes using airborne remote sensing. *ISPRS-J. Photogramm. Remote Sens.* 109,
651 165-177.

652 O'Connell, J., Connolly, J., Vermote, E.F., Holden, N.M., 2013. Radiometric normalization for change
653 detection in peatlands: a modified temporal invariant cluster approach. *Int. J. Remote Sens.* 34,
654 2905-2924.

655 Pena, J.M., Gutierrez, P.A., Hervas-Martinez, C., Six, J., Plant, R.E., Lopez-Granados, F., 2014. Object-
656 Based Image Classification of Summer Crops with Machine Learning Methods. *Remote Sens.* 6,
657 5019-5041.

658 R Core Team, 2016. R: A language and environment for statistical computing. R Foundation for
659 Statistical Computing, Vienna, Austria, <http://www.R-project.org>.

660 Silla, C.N., Freitas, A.A., 2011. A survey of hierarchical classification across different application
661 domains. *Data Min. Knowl. Discov.* 22, 31-72.

662 Tarquini, S., Isola, I., Favalli, M., Mazzarini, F., Bisson, M., Pareschi, M.T., Boschi, E., 2007.
663 TINITALY/01: a new Triangular Irregular Network of Italy.

664 Tarquini, S., Vinci, S., Favalli, M., Doumaz, F., Fornaciai, A., Nannipieri, L., 2012. Release of a 10-m-
665 resolution DEM for the Italian territory: Comparison with global-coverage DEMs and anaglyph-
666 mode exploration via the web. *Computers & Geosciences* 38, 168-170.

667 Thoonen, G., Spanhove, T., Vanden Borre, J., Scheunders, P., 2013. Classification of heathland
668 vegetation in a hierarchical contextual framework. *Int. J. Remote Sens.* 34, 96-111.

669 Thuiller, W., Araujo, M.B., Lavorel, S., 2004. Do we need land-cover data to model species
670 distributions in Europe? *Journal of Biogeography* 31, 353-361.

671 Tomaselli, V., Adamo, M., Veronico, G., Sciandrello, S., Tarantino, C., Dimopoulos, P., Medagli, P.,
672 Nagendra, H., Blonda, P., 2016. Definition and application of expert knowledge on vegetation
673 pattern, phenology, and seasonality for habitat mapping, as exemplified in a Mediterranean
674 coastal site. *Plant Biosystems - An International Journal Dealing with all Aspects of Plant Biology*,
675 1-13.

676 Tomaselli, V., Dimopoulos, P., Marangi, C., Kallimanis, A.S., Adamo, M., Tarantino, C., Panitsa, M.,
677 Terzi, M., Veronico, G., Lovergine, F., Nagendra, H., Lucas, R., Mairota, P., Mucher, C.A., Blonda,
678 P., 2013. Translating land cover/land use classifications to habitat taxonomies for landscape
679 monitoring: a Mediterranean assessment. *Landscape Ecology* 28, 905-930.

680 Tso, B., Mather, P., 2009. Classification methods for remotely sensed data. CRC Press, Taylor &
681 Francis Group.
682 Xie, Y., Sha, Z., Yu, M., 2008. Remote sensing imagery in vegetation mapping: a review. J. Plant Ecol.
683 1, 9-23.

684

685 SUPPORTING INFORMATION

- 686 • Supporting Information S1 – The image processing procedure, the derived indices and the
687 environmental explanatory variables.
- 688 • Supporting Information S2 – Information on the Hierarchical F measure
- 689 • Supporting Information S3 –Results from the nine mixed effect models.
- 690 • Supporting information S4 – excel file with example confusion matrices and frequency of classes
691 in the training and validation sets.
- 692 • Supporting information S5 – zip file with the '*HieRanFor*' package, verified to work on R version
693 3.3.1 (July 2017).

694

695 TABLES

696 **Table 1:** The top 5 pairs of habitat/land-cover that were responsible for the highest number of classification
 697 errors in the validation sets in each of the three classification schemes. Results for the flat, hierarchical
 698 multiplicative (H.Mult) and hierarchical stepwise (H.Step) are given as the mean over all 15 runs. Number of
 699 errors is the total number of mix-up (i.e., number of times A was classified as B plus the number of times B was
 700 classified as A).

Habitat/Land-cover 1	Habitat/Land-cover 2	Number of errors (rank)			
		Flat	H.Mult	H.Step	
CORINE		Total errors:	3907	3956	4239
Coniferous forest	Road and rail networks and associated land	1718 (1)	1890 (1)	1964 (1)	
Olive groves	Permanently irrigated land	438 (2)	488 (2)	508 (2)	
Olive groves	Road and rail networks and associated land	412 (3)	400 (3)	389 (3)	
Coniferous forest	Discontinuous urban fabric	312 (4)	181 (6)	299 (5)	
Coniferous forest	Olive groves	245 (5)	253 (5)	250 (6)	
Coniferous forest	Permanently irrigated land	225 (6)	265 (4)	318 (4)	
FAO-LCCS		Total errors:	5107	5310	4857
Plantations: needleleaved evergreen tree crops- monoculture + rainfed	Paved roads	1415 (1)	1448 (1)	877 (2)	
Plantations: needleleaved evergreen tree crops- monoculture + rainfed	Orchards: broadleaved evergreen tree crops- monoculture + rainfed	1144 (2)	1213 (2)	1469 (1)	
Orchards: broadleaved evergreen tree crops- monoculture + rainfed	Paved roads	773 (3)	752 (3)	537 (4)	
Orchards: broadleaved evergreen tree crops- monoculture + rainfed	Fields of irrigated no graminoid crops + one additional crop	500 (4)	544 (4)	602 (3)	
Orchards: broadleaved evergreen tree crops- monoculture + rainfed	Scattered industrial or other areas	312 (5)	294 (6)	271 (7)	
Temporarily flooded land with perennial closed tall grasslands	Temporarily flooded land with <i>Aphyllous</i> closed dwarf shrubs	299 (6)	369 (5)	369 (5)	
EUNIS		Total errors:	7025	6885	6955
Fen <i>Cladium mariscus</i> beds	Marine saline beds of <i>Phragmites australis</i>	3177 (1)	3115 (1)	3129 (1)	
Native conifer plantations	Road networks	1711 (2)	1635 (2)	1635 (2)	
<i>Olea europaea</i> groves	Road networks	467 (3)	425 (4)	442 (3)	
<i>Olea europaea</i> groves	Arable land with unmixed crops grown by low-intensity agricultural methods	421 (4)	431 (3)	389 (5)	
Native conifer plantations	Scattered residential buildings	316 (5)	208 (7)	278 (6)	
Native conifer plantations	<i>Olea europaea</i> groves	235 (6)	362 (5)	407 (4)	

701

702

703

704 FIGURE LEGENDS

705 **Figure 1: Main concept of hierarchical random forest.**

706 Flow chart for creating: (A) a single classification tree, (B) a single randomForest model, and (C) a single
707 hierarchical randomForest model. See text for details. Prop' – proportion; OoB-- Out-of-Bag.

708 **Figure 2: The results for CORINE.**

709 The top panel is an example of the observed map and the maps predicted by a flat RF model, by Hierarchical
710 multiplicative RF model and by hierarchical stepwise RF. Land-covers legend is given in the class hierarchy
711 (lower right panel) alongside the location of seven local classifiers (e.g., C.1). The lower left panel shows the
712 variable importance values in each local classifier and in the flat model. The variables are divided to
713 environmental variables (blue), variables derived from the remote sensing images (black) and the raw
714 reflectance values of the two remotely-sensed images (red).

715 **Figure 3: The results for FAO-LCCS.**

716 The top panel is an example of the observed map and the maps predicted by a flat RF model, by Hierarchical
717 multiplicative RF model and by hierarchical stepwise RF. Land-covers legend is given in the class hierarchy
718 (lower right panel) alongside the location of eight local classifiers (e.g., C.1). The lower left panel shows the
719 variable importance values in each local classifier and in the flat model. The variables are divided to
720 environmental variables (blue), variables derived from the remote sensing images (black) and the raw
721 reflectance values of the two remotely-sensed images (red).

722 **Figure 4: The results for EUNIS.**

723 The top panel is an example of the observed map and the maps predicted by a flat RF model, by Hierarchical
724 multiplicative RF model and by hierarchical stepwise RF. Habitats legend is given in the class hierarchy (lower
725 right panel) alongside the location of eleven local classifiers (e.g., C.1). The lower left panel shows the variable
726 importance values in each local classifier and in the flat model. The variables are divided to environmental
727 variables (blue), variables derived from the remote sensing images (black) and the raw reflectance values of
728 the two remotely-sensed images (red).

729 **Figure 5: Summary of models' performance.**

730 Performance of the FRF (Flat), HRF with multiplicative majority rule (H.Mult) and HRF with stepwise majority
731 rule (H.Step) in each classification scheme, for the training and validation sets. Box plots represent the 25, 50
732 and 75% percentiles of the 15 runs, while the whiskers give the 1.5 IQR. Outliers are given as points while the
733 mean is given as a triangle.

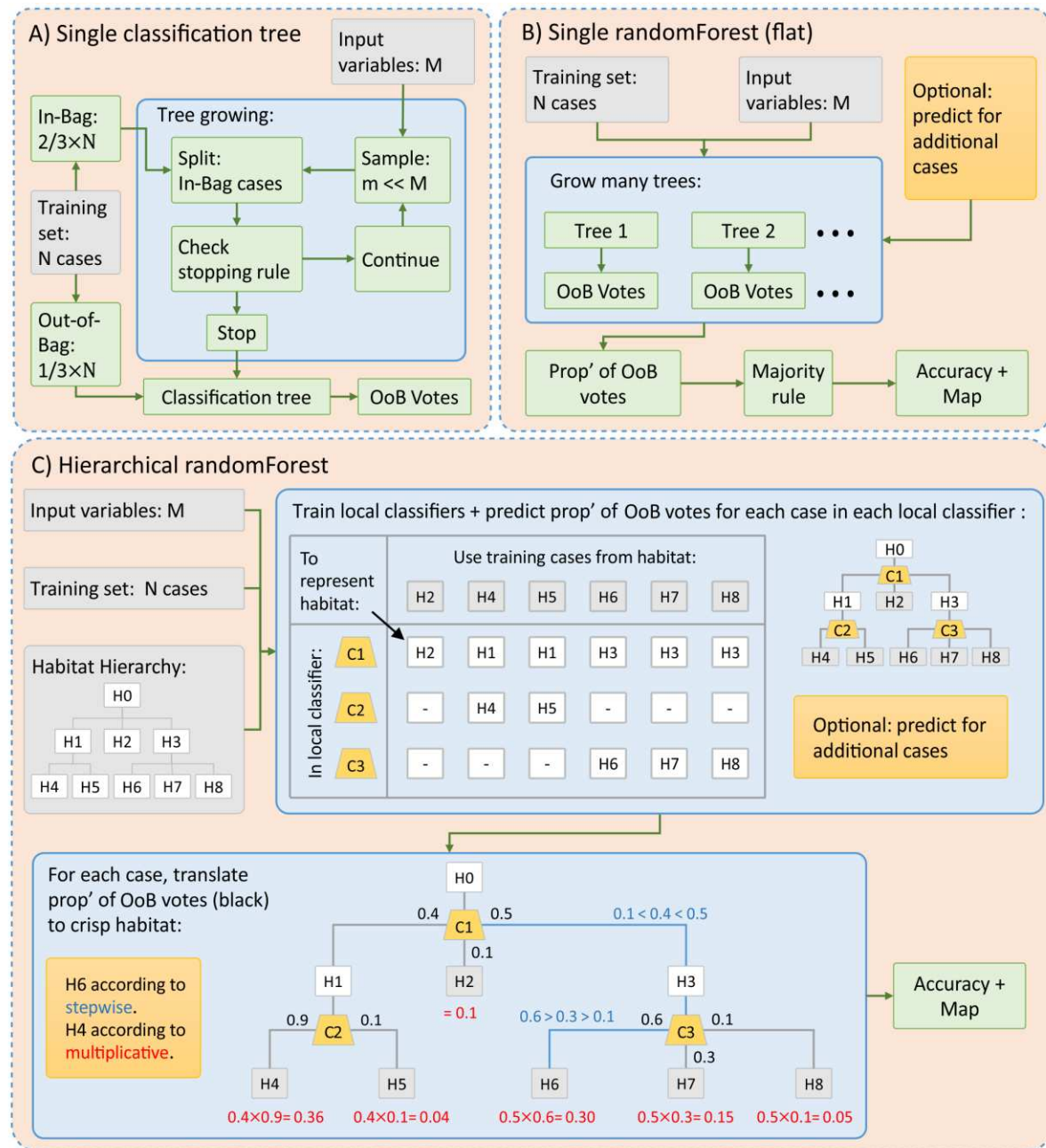
734

735

736

737

738 Figure 1:



739

740

741

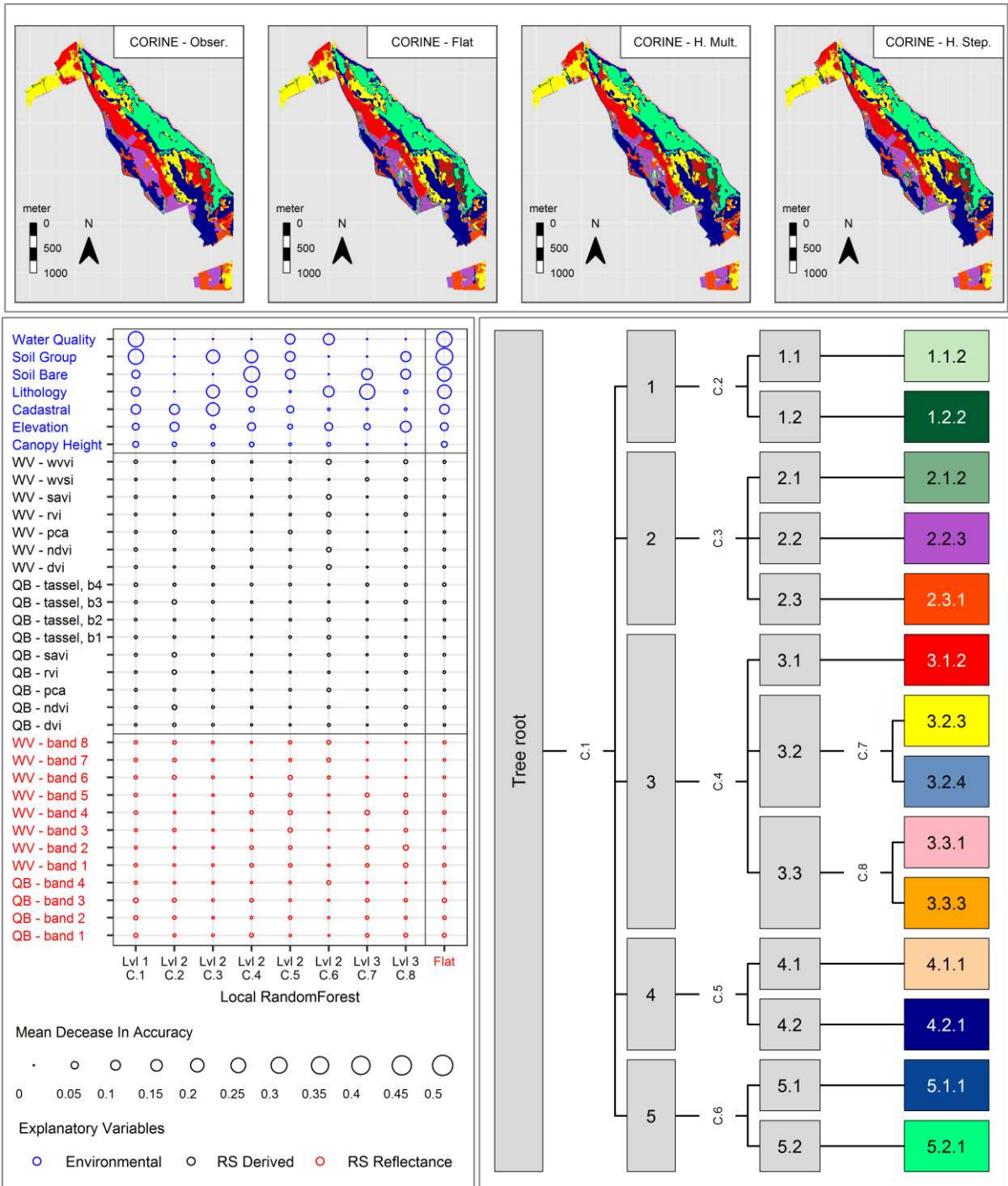
742

743

744

745

746 Figure 2:



747

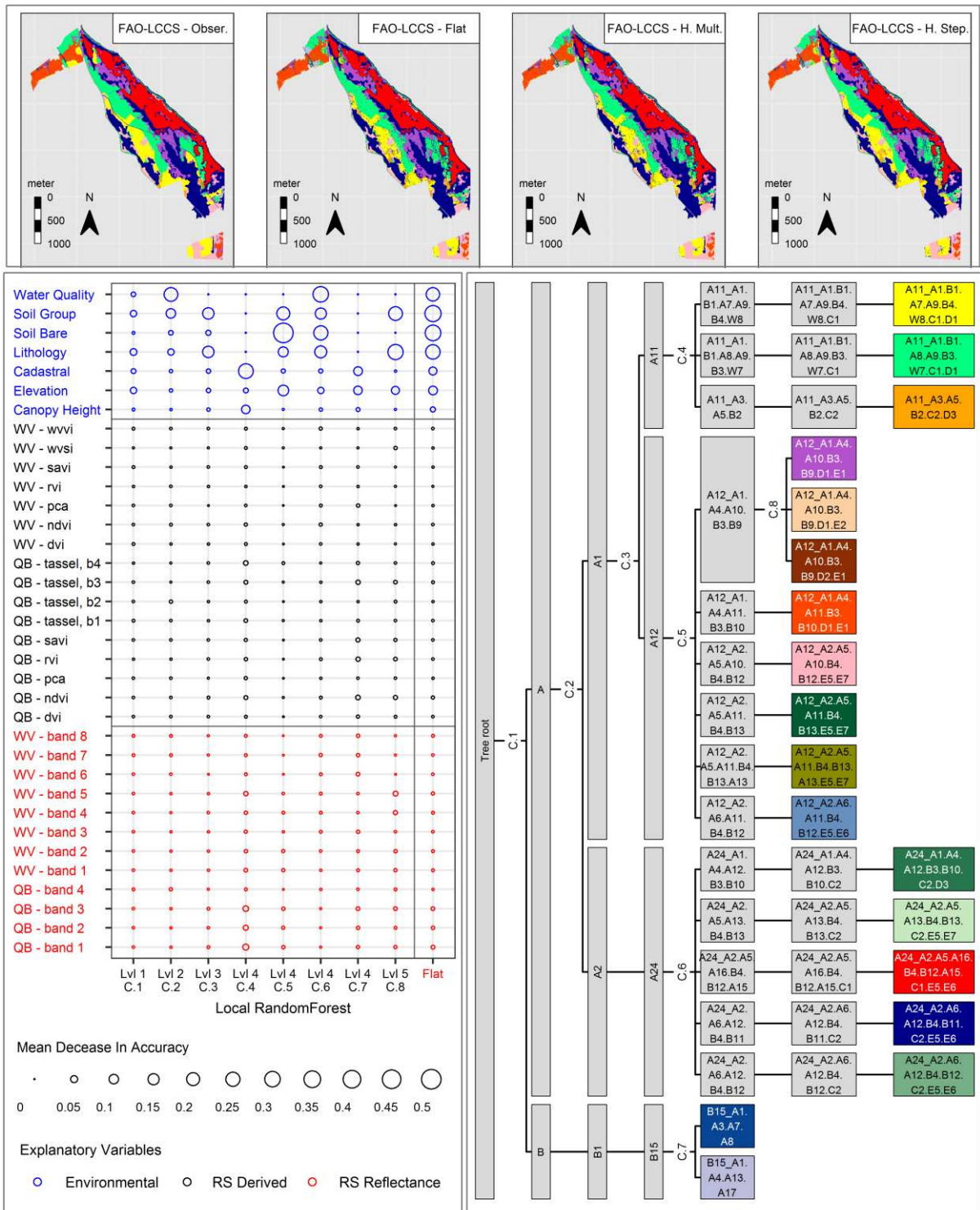
748

749

750

751

752 Figure 3:



753

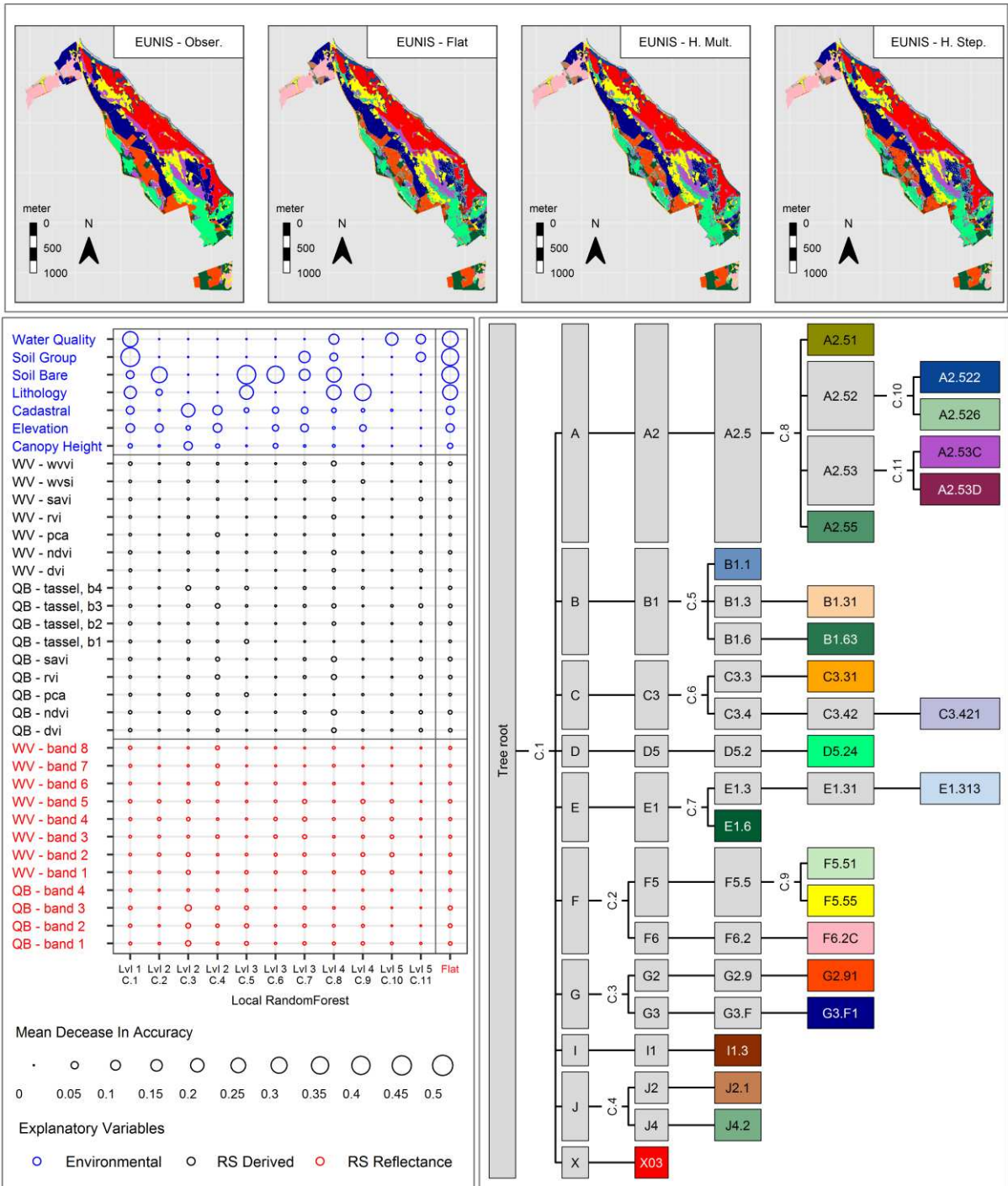
754

755

756

757

758 Figure 4:



759

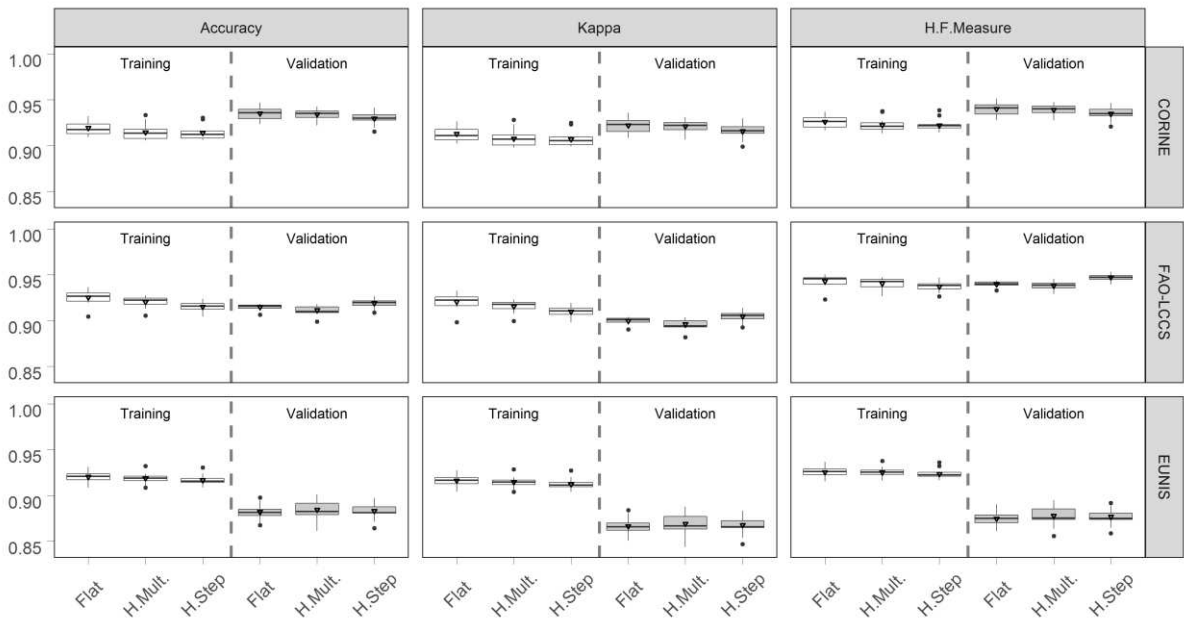
760

761

762

763

764 Figure 5:



765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

780

