

1 **Brief heading: Chemocoding: identification tool for species delimitation in species-**  
2 **rich genera in the tropics.**

3  
4  
5 **Title: Chemocoding as an identification tool where morphological- and DNA-based**  
6 **methods fall short: *Inga* as a case study**

7  
8 **María-José Endara<sup>1,2</sup>, Phyllis D. Coley<sup>1,3</sup>, Natasha L. Wiggins<sup>4</sup>, Dale L. Forrister<sup>1</sup>,**  
9 **Gordon C. Younkin<sup>1</sup>, James A. Nicholls<sup>5</sup>, R. Toby Pennington<sup>6</sup>, Kyle G. Dexter<sup>6,7</sup>,**  
10 **Catherine A. Kidner<sup>5,6</sup>, Graham N. Stone<sup>5</sup>, and Thomas A. Kursar<sup>1,3</sup>**

11  
12 <sup>1</sup>Department of Biology, University of Utah, Salt Lake City, UT 84112-0840, USA

13 <sup>2</sup>Centro de Investigación de la Biodiversidad y Cambio Climático (BioCamb) e  
14 Ingeniería en Biodiversidad y Recursos Genéticos. Facultad de Ciencias de Medio  
15 Ambiente, Universidad Tecnológica Indoamérica, Quito EC170103, Ecuador

16 <sup>3</sup>Smithsonian Tropical Research Institute Box 0843-03092, Balboa, Ancón, Republic of  
17 Panamá

18 <sup>4</sup>School of Biological Sciences, University of Tasmania, Sandy Bay TAS 7001, Australia

19 <sup>5</sup>Ashworth Labs, Institute of Evolutionary Biology, School of Biological Sciences,  
20 University of Edinburgh, Edinburgh EH9 3JY, UK

21 <sup>6</sup>Royal Botanic Garden Edinburgh, Edinburgh EH3 5LR, UK

22 <sup>7</sup>School of GeoSciences, University of Edinburgh, Edinburgh EH9 3FF, UK

23  
24 **Correspondence:** María-José Endara; Department of Biology, University of Utah, Salt  
25 Lake City, UT 84112-0840, USA; [majo.endara@utah.edu](mailto:majo.endara@utah.edu), tel. +01 (801) 581-7086

26  
27 **Word count:** Main body: 5,704; Introduction: 897; Materials & Methods: 1,480; Results:  
28 1,245; Discussion: 1,945; Acknowledgments: 82.

29 **Tables:** 1

30 **Figures:** 5 (all color)

31 **Supplementary information:** Tables: 3; Figures: 9

32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62

## Summary

- The need for species identification and taxonomic discovery has led to the development of innovative technologies for large-scale plant identification. DNA barcoding has been useful, but fails to distinguish among many species in species-rich plant genera, particularly in tropical regions. Here, we show that chemical fingerprinting, or “chemocoding”, has great potential for plant identification in challenging tropical biomes.
- Using untargeted metabolomics in combination with multivariate analysis, we constructed species-level fingerprints, which we define as chemocoding. We evaluated the utility of chemocoding with species that were defined morphologically and subject to next-generation DNA sequencing in the diverse and recently radiated neotropical genus, *Inga* (Leguminosae), both at single study sites and across broad geographic scales.
- Our results show that chemocoding is a robust method for distinguishing morphologically similar species at a single site and for identifying widespread species across continental-scale ranges.
- Given that species are the fundamental unit of analysis for conservation and biodiversity research, the development of accurate identification methods is essential. We suggest that chemocoding will be a valuable additional source of data for a quick identification of plants, especially for groups where other methods fall short.

**Key-words:** Chemocoding, *Inga*, metabolomics, species identification, tropical forests

63

64 **Introduction**

65 Cataloguing the world's plant diversity has been a challenge for centuries, and because of  
66 accelerated anthropogenic extinctions, the rapid documentation of biodiversity is more  
67 critical than ever (Mace, 2004; Valentini *et al.*, 2009; Vernooy *et al.*, 2010; Cristescu,  
68 2014). This is particularly true in tropical rainforests, where the high diversity and co-  
69 occurrence of morphologically and ecologically similar congeneric species have  
70 presented significant challenges for identification (Gonzalez *et al.*, 2009; Dexter *et al.*,  
71 2010; Kress *et al.*, 2015; Liu *et al.*, 2015). Plant species have typically been identified by  
72 botanical experts based on morphological characteristics, or more recently, by sequencing  
73 several chloroplast DNA regions or the internal transcribed spacer of nuclear ribosomal  
74 DNA (ITS), referred to as "DNA barcoding" (Gemeinholzer *et al.*, 2006; Hollingsworth  
75 *et al.*, 2009; Chen *et al.*, 2010; Kress *et al.*, 2010; China Plant BOL Group *et al.*, 2011;  
76 Hollingsworth *et al.*, 2011). Here we present evidence that chemical fingerprinting or  
77 "chemocoding" can be another tool for species identification in confusing and/or closely  
78 related tree species in challenging, hyperdiverse biomes.

79 No single identification method is without drawbacks. Morphological methods are  
80 often labor-intensive, rely upon taxonomic expertise and are the most prone to subjective  
81 errors, particularly where phenotypic plasticity and cryptic taxa are prevalent (de  
82 Carvalho *et al.*, 2005; Costion *et al.*, 2011). Although DNA barcoding is rapid and  
83 straightforward, it can fail to distinguish closely related plant species because of  
84 insufficient sequence divergence in standard barcode markers (Kress *et al.*, 2009; Dexter  
85 *et al.*, 2010, Liu *et al.*, 2015) and may lead to faulty identifications in genera where  
86 species are of recent origin (Razafimandimbison *et al.*, 2004, Naciri & Linder, 2015;  
87 Pennington & Lavin, 2015). These limitations are often exacerbated in highly diverse  
88 tropical systems, as the proportion of species belonging to young, species-rich genera is  
89 high (Richardson *et al.*, 2001). For example, in a subtropical Chinese forest, for the 44%  
90 of species belonging to genera with more than two species, >50% shared barcoding  
91 sequences and could not be distinguished (Liu *et al.*, 2015). This led to an overall species  
92 resolution of only 67%. A similar problem is encountered in New World tropical  
93 rainforests, where DNA barcoding cannot reliably discriminate species within

94 ecologically important, species-rich genera such as *Inga*, *Ficus* and *Piper* (Gonzalez *et*  
95 *al.*, 2009; Kress *et al.*, 2009). While standard DNA barcoding for problematic groups can  
96 be improved by adding data for additional loci, recently diverged species will always be  
97 hard to distinguish using sequence data and no standardized marker sets for such  
98 extended DNA barcoding exist.

99 Many studies assessing diversity in surveys or plots for conservation or basic  
100 science rely on identifying all individuals in a plot to species-level, including the large  
101 majority of individuals that are without flowers or fruits (Dexter *et al.*, 2010). Plots in the  
102 tropics represent a daunting task, and are still faced with problematic identifications  
103 despite extremely well-trained botanists. A related problem is the difficulty of achieving  
104 uniform species identifications across multiple sites. For example, in a recent, extensive  
105 analysis of the identifications for eight genera, the three genera with the highest error  
106 rates were *Andira* and *Tachigali* (approximately 50%) and *Inga* (about 40%; Baker *et al.*,  
107 2017). And yet another substantial challenge is correlating the identities to species level  
108 for saplings, small trees and adult trees, because ontogenetic changes in leaf morphology  
109 may be considerable and juveniles do not bear flowers or fruits. For example, a  
110 consortium to understand forest dynamics has established 63 plots around the world  
111 where all woody plants >1 cm DBH are mapped and identified, the majority of which are  
112 juveniles ([www.forestgeo.si.edu](http://www.forestgeo.si.edu)). Thus, the issues and errors associated with  
113 morphological species identifications of thousands of trees in the tropical forests are  
114 serious.

115 Given that species are a fundamental unit of analysis for conservation, for  
116 quantifying biodiversity, and for understanding ecological and evolutionary processes,  
117 the development of accurate methods for identifying them is essential. In this paper, we  
118 suggest that chemical fingerprinting (here termed chemocoding) can provide an  
119 additional identification tool for species identification in a species-rich tropical tree  
120 genus, particularly for morphologically confusing or cryptic species. We examine its  
121 utility both for distinguishing species within a single site, and for characterizing within-  
122 species variation over wider geographic scales. Moreover, chemocoding may be  
123 inexpensive enough to allow for every individual tree to be tested.

124 We test the potential of chemocoding for species identification within *Inga* Mill.  
125 (Leguminosae, Mimosoideae) because species in this genus are difficult to distinguish  
126 morphologically and show insufficient variation in barcoding sequences (Richardson *et*  
127 *al.*, 2001; Kress *et al.*, 2009; Dexter *et al.*, 2010; Dick & Webb, 2012). *Inga* is one of the  
128 most abundant and diverse Neotropical genera in lowland forest communities (Valencia  
129 *et al.*, 1994; ter Steege *et al.*, 2013), is widely distributed, and has undergone recent, rapid  
130 diversification (Richardson *et al.*, 2001). For *Inga*, genetic and morphological  
131 differentiation of closely related species is low and the identification of a species can  
132 therefore be difficult.

133 We propose the use of small, defense-related chemical markers characterized via  
134 untargeted metabolomics in combination with multivariate analysis for the construction  
135 of a phytochemical, species-level fingerprint, which we define as chemocoding. We  
136 evaluate the units defined by chemocoding with those defined morphologically in a  
137 recent taxonomic monograph (Pennington, 1997) as well as with those defined using  
138 next-generation DNA sequencing data of many hundreds of nuclear genes (Nicholls *et*  
139 *al.*, 2015 & unpubl. data).

140

## 141 **Material and Methods**

### 142 **Study sites**

143 Samples were collected at five sites that include a wide range of soils but very similar  
144 climates throughout the Amazon and Panama (Fig. 1). Barro Colorado Island is a field  
145 station administered by the Smithsonian Research Tropical Institute located in the  
146 Panama Canal (9°N 80°W). It is a lowland moist forest with 2649 mm of precipitation a  
147 year and 4-month dry season with mean monthly temperatures of 27°C (Leigh, 1999).  
148 The other four sites do not have a pronounced dry season. The Nouragues Ecological  
149 Research Station, French Guiana (4°N 53°W) is located inside the Nouragues National  
150 Reserve on the Guiana Shield. The mean annual precipitation is 2990 mm and the mean  
151 annual temperature is 26.3°C (Grimaldi & Riera, 2001). Tiputini Biodiversity Station is  
152 located in the eastern lowland Ecuadorian Amazon (0°S 75W), inside the Yasuni  
153 Biosphere Reserve. The climate is humid and aseasonal, with an annual precipitation of  
154 3320 mm and an average annual temperature of 26° C (Valencia *et al.*, 2004). Kilometer

155 41 (KM41, 2°S 59°W) is a field station of the Biological Dynamics of Forest Fragments  
156 project located near Manaus, Brazil. The mean annual temperature is 26°C and average  
157 annual precipitation is 2651 mm (Radtke *et al.*, 2007). Los Amigos Biological Station is  
158 located in the southeastern lowland Peruvian Amazon, in the Madre de Dios Department  
159 (13°S 70°W). The mean annual rainfall is between 2700 to 3000 mm. Due to winter cold  
160 spells, the daily minimum temperature can drop to less than 10°C; the mean monthly  
161 temperature range is from 21 to 26°C (Pitman, 2007). For simplicity in the text, each site  
162 will be referred to by the country only.

163

### 164 **Study species**

165 We examined saplings of *Inga* because this size class is the most frequently censused for  
166 ecological research and also is the size class that can be very difficult to identify.  
167 Morphological identifications were based on the most recent taxonomic *Inga* monograph  
168 (Pennington, 1997), and made by four researchers (MJE, TAK, PDC and KGD) who  
169 have worked in the field identifying *Inga* for about five decades collectively. They  
170 consulted with the botanists working in the 50-ha CTFS (Center for Tropical Forest  
171 Science) plots in Panama and Ecuador and the plots in Nouragues, French Guiana.

172

### 173 **Sampling**

174 We determined the power of chemocoding for discriminating among species and among  
175 geographically disjunct populations within species. We included different taxa that are  
176 similar in vegetative morphology and are very difficult to identify in the absence of  
177 reproductive structures. Our examples include cases in which these coexist at the same  
178 site or in two well separated sites. We also included cases of one species that is  
179 morphologically uniform in collections from up to four sites (Table 1).

180 In addition, we tested the ability of chemocoding to correctly distinguish simultaneously  
181 between populations of several dozens of species at a regional scale (see Random Forest  
182 Analysis, below). More examples can be found in the supplementary information.

183

### 184 **Collections**

185 For each species, samples of expanding leaves were collected from five saplings, 0.5 – 4  
186 m in height, in the shaded understory. We focused on expanding leaves as part of a study  
187 of plant-herbivore interactions and also because secondary metabolites are at greater  
188 concentration during the expansion stage than in leaves that have matured and toughened  
189 (Wiggins *et al.*, 2016). For each sapling, we collected leaves that were between 20% and  
190 80% of the average maximum size. Fresh leaves were dried at room temperature with  
191 fans and silica gel during 24-48h, transported to the University of Utah, and stored at -  
192 20°C. For DNA analysis we typically included one sample per species per site.

193

### 194 **Metabolomic analysis**

195 Metabolites were extracted and analyzed following the protocol of Wiggins *et al.*, (2016),  
196 specifically designed for secondary metabolites having intermediate polarity. In *Inga*,  
197 these are mainly phenolics and saponins. Briefly, 100 mg of ground leaves were extracted  
198 in 1.0 ml of extraction buffer (44.4 mM ammonium acetate (pH 4.8): acetonitrile, 60:40,  
199 v/v). After extraction for 5 min and centrifugation (13 793 x g) for 5 min, the supernatant  
200 was transferred to a glass vial and the extraction repeated. The extracts were diluted  
201 fivefold by combining 200 µL of crude extract with 790 µL of acetonitrile:water (60/40,  
202 v/v) plus 10 µL of internal standard (1mg mL<sup>-1</sup> biochanin A in acetonitrile:water, 50/50).  
203 Soluble metabolites were analyzed by ultra-performance liquid chromatography coupled  
204 to mass spectrometry (UPLC-MS) using an Acquity UPLC<sup>®</sup> *I-Class* system and a Xevo<sup>®</sup>  
205 G2 Q-ToF MS equipped with LockSpray and an electrospray ionization source (Waters,  
206 Milford, MA). Data were collected in negative ionization mode.

207 Raw data from the UPLC-MS were processed for peak detection, peak alignment  
208 and peak filtering using MassLynx (Waters) and the R package XCMS (Smith *et al.*,  
209 2006; Tautenhahn *et al.*, 2008; Benton *et al.*, 2010). The parameters used were: peak  
210 detection method “*centWave*” (ppm=15, peakwidth=c(5,12), snthresh=5); peak grouping  
211 method “*density*” (bw=2); retention time correction method “*obiwarp*”; integrate areas  
212 of missing peaks method “*chrom*”. XCMS processing was performed for each species  
213 independently, with five leaf samples included as replicates. The results obtained by  
214 XCMS were post-processed in the R package CAMERA to assign the various ions  
215 derived from one compound (termed “features”) to that compound (Kuhl *et al.*, 2012).

216 This uses a defined set of rules for linking the precursor ion with adducts and neutral  
217 losses (Table S1 for a list of these). The parameters used were: peak grouping after  
218 retention time “*groupFWHM*” (perfwhm: 0.8); verify grouping “*groupCorr*”; annotate  
219 isotopes “*findIsotopes*”; annotate adducts “*findAdducts*” (polarity= “negative”). For each  
220 case study, the resulting peak tables for each species were combined into a single peak  
221 table (*m/z* and retention time for each peak) using the R package metaXCMS  
222 (Tautenhahn *et al.*, 2011) with the following parameters: peak filtering: none; *m/z* and  
223 retention time tolerance: 0.05 and 12 seconds respectively. Peak tables are stored in  
224 MetaboLights (<https://www.ebi.ac.uk/metabolights/>), a publicly available database.

225

## 226 **Statistical Analysis**

227 The variation in metabolites across samples was quantified using unsupervised  
228 multivariate methods (no prior classification of samples), which is suitable for  
229 metabolomics data. For the first four case studies, we chose methods for data reduction  
230 and pattern recognition that group and visualize samples according to their similarities  
231 without prior assignment of samples to classes (Bartel *et al.*, 2013).

232 In order to visualize grouping patterns across samples, we used hierarchical  
233 clustering with multiple agglomerative algorithms because this method works well for a  
234 limited number of species (classes) and provides statistical power (Embrechts *et al.*,  
235 2013). Peak intensities, or the total ion current (TIC), were normalized by dividing by the  
236 sum of the TIC for all features in the chromatogram of a sample. Subsequently, we fitted  
237 a hierarchical clustering model to the normalized data with 10000 permutations using the  
238 R package pvclust (Suzuki & Shimodaira, 2014). The hierarchical clustering was  
239 performed using the Pearson’s correlation similarity measure, a routine method adopted  
240 for “omics” data (Reeb *et al.*, 2015). The clustering algorithm selection for each analysis  
241 was based on the correlation between the original distance matrix and the patristic  
242 distance in the hierarchical cluster diagram. Clusters with AU (Approximately Unbiased)  
243 p-values of  $\geq 95\%$  are considered to be strongly supported by the data. For more details  
244 see Wiggins *et al.*, (2016).

245 In addition, to test the overall accuracy of chemocoding to identify samples when  
246 presented with a very large number of species (classes) simultaneously, we used

247 supervised statistical learning methods. Specifically, we chose Random Forest Analysis,  
248 which is a powerful classification method for multivariate datasets with many weak  
249 predictor variables along with a large number of species (cases). This method has been  
250 widely adopted in remote sensing, high dimensional biological data (various “omics), and  
251 ecology (Breiman, 2001; Lawrence *et al.*, 2006; Cutler *et al.*, 2007; Cutler *et al.*, 2009).  
252 Based on models that we constructed with the metabolomics data, we used Random  
253 Forest in order to predict how well samples can be classified to species. For this, a single  
254 sample-by-compound matrix was generated using XCMS as described above. A total of  
255 1000 trees was generated for each Random Forest Model and 100 variables were used at  
256 each split, which was sufficient to arrive at a model with minimal prediction error  
257 (Breiman, 2011). We performed the analyses for 82 species with samples selected from a  
258 single site, as well as for 26 species found at two to four sites. Analyses were performed  
259 using the randomForest R package (Liaw & Weiner, 2002). R code for all of the analyses  
260 is provided in the Supplementary Section (Methods S1-S2).

261

## 262 **Next-generation DNA sequence data**

263 To determine the accuracy of our approach, we compared delimitations based on  
264 chemocoding with the first resolved phylogeny of *Inga*, accomplished through targeted  
265 enrichment and sequencing of 194 loci (259,313 bases; Nicholls *et al.*, 2015 & unpubl.  
266 data). Due to *Inga*'s recent, rapid radiation (Richardson *et al.*, 2001), a previous  
267 phylogeny with over 6kb of plastid and nuclear DNA sequence did not resolve species-  
268 level relationships fully (Kursar *et al.*, 2009).

269

## 270 **Results**

271 *Case study 1: Two morphologically confusing or cryptic species that are present at the*  
272 *same site*

273 *Inga alata* Benoist and *Inga pezizifera* Benth are sister species (Nicholls *et al.*, 2015, &  
274 unpubl. data; Fig. 2c). The saplings can only be successfully differentiated by expert field  
275 workers using subtle differences in the shape of the extrafloral nectaries, the number of  
276 leaflets, the number of primary lateral leaf veins, and the color of the expanding leaves  
277 (Fig. S1). In addition, they differ in the morphology of their inflorescences (Pennington,

278 1997), but this feature is not available in the saplings studied by many ecologists. We  
279 investigated how chemocoding might be useful to separate these two confusing species in  
280 Nouragues, French Guiana, where one species is often found meters away from the other.

281 Consistent with DNA sequence differences (Fig. 2c), chemocoding accurately  
282 determined species limits for five saplings each of the two species. The profiles of  
283 secondary metabolites showed visually evident differences between species (Fig. 2a).  
284 Hierarchical clustering of UPLC-MS metabolomics data (98% AUP value, Fig. 2b)  
285 clustered the samples into two distinct groups, one for each species.

286 We evaluated four further groups of species that are hard to separate  
287 morphologically, and coexist at a single site (available in the supplementary section). For  
288 three of these, chemocoding-based separation agreed with DNA sequence differences  
289 (Fig. S2: *I. coruscans* & *I. laurina* in Peru, Fig. S3: *I. umbellifera* T50, T72 and T73 in  
290 Ecuador (where T numbers correspond to codes for morphotypes in Tiputini, Ecuador),  
291 and Fig. S4: *I. chartaceae* & *I. sapindoides* in Ecuador). For the fourth supplemental  
292 case, chemocoding found substantial differences in secondary metabolites between two  
293 morphotypes of *I. leiocalycina* (T65 and T86, Fig. S5). T86 occurs in terra firme forests  
294 and T65 in floodplains in the Ecuadorian Amazon, whereas DNA data placed these two  
295 morphotypes into a monophyletic group (Fig. S5d). This may be the result of plasticity in  
296 response to differences in habitat (although we found no intermediate morpho- or  
297 chemotypes). Alternatively, these may reflect strong selection in the face of gene flow  
298 across this environmental gradient or these may be distinct species.

299

300 *Case study 2: A single species that shows morphological variation within a site*

301 *Inga acreana* Harms is a widely-distributed species across South America, from the  
302 Guyanas to the Amazon Basin in Colombia, Ecuador, Peru and Bolivia (Pennington,  
303 1997; Pennington & Revelo, 1997). In the Tiputini Biological Station in Ecuador, it  
304 comprises two distinct types that have the same morphology, but that differ subtly in the  
305 color of the expanding leaves (T28 and T56; T numbers are codes for morphotypes; Fig.  
306 S1), and co-occur in floodplain habitats. We used chemocoding of five saplings of each  
307 of these two leaf variants to determine if they were the same or different chemotypes.

308 The two morphotypes showed no consistent metabolomic differences (Fig. 3a).  
309 Hierarchical clustering models fitted to the UPLC-MS data reveal no separate clusters  
310 (Fig. 3b). These results agree with DNA data; the phylogeny from targeted enrichment  
311 data placed these accessions representing these two morphotypes into the same, otherwise  
312 unstructured, monophyletic group (Fig. 3c). Most likely, the color difference between the  
313 two morphotypes is due to a difference in anthocyanin production, a form of intraspecific  
314 variation sometimes seen in expanding leaves.

315

316 *Case study 3: Distinguishing among two morphologically similar species across three*  
317 *sites*

318 Cross-checking identifications amongst morphologically similar tree species that occur at  
319 different sites is particularly challenging in tropical forests (Baker *et al.*, 2017). *Inga*  
320 *brachystachys* Ducke and *Inga obidensis* Ducke are closely related (Fig. 4d), with  
321 overlapping distributions across the Amazon (Pennington, 1997; Pennington & Revelo,  
322 1997). They are morphologically very similar in the vegetative state, making it  
323 problematic to assign accurate species names (Fig. S1).

324 Chemocoding and hierarchical clustering of five saplings from each of the three  
325 sites delimited the samples into two groups, separating the French Guiana and Brazil  
326 samples from those collected in Peru (100% AUP value, Fig. 4c). Together, DNA  
327 sequence data (Fig. 4c), chemocoding and morphology suggest that samples collected in  
328 Brazil and French Guiana are a single species, *I. obidensis*, and that the chemically  
329 distinct Peruvian samples may represent a different, as yet unidentified species that, in its  
330 vegetative morphology, is similar to *I. brachystachys*.

331

332 *Case study 4: Identification of a widespread species across its range*

333 In order to assess the variation of chemocoding profiles across geographic space in a  
334 widespread species, we collected data on five saplings per population of *Inga auristellae*  
335 Harms, a species that occurs across northern South America. These came from four  
336 geographically separated populations: French Guiana, Ecuador, Brazil, and Peru (Fig. 1).

337 The different populations showed consistency in their chemistry across their wide  
338 geographic range (Fig. 5a). The hierarchical clustering model shows no significant

339 differences between *I. auristellae* populations from different geographic areas (94% AUP  
340 value, Fig. 5b). The DNA sequences analyses show that although the four populations  
341 belong to the same species, there is a strong geographic structure, with an initial split into  
342 east vs west Amazonia, and then within each of these groups the samples cluster by  
343 population (Fig. 5c). Thus, chemocoding shows consistency in species characterization,  
344 despite some evidence of genetic population structure without corresponding structuring  
345 of chemistry. A similar pattern was observed for other widespread species: *I. pezizifera*  
346 (Fig. S6). A second widespread species, *I. alba*, showed no geographic structure either in  
347 chemistry or DNA (Fig. S7).

348

349 *Case 5: Identification of a large number of species, including comparisons among*  
350 *populations of widespread species*

351 The detailed case studies presented above allow a small number of samples to be grouped  
352 by chemical similarity and statistically validated without prior classification of the  
353 samples (unbiased). Additionally, we sought to evaluate how often chemocoding could  
354 correctly identify samples to species when we include a large number of classes (species).  
355 To this end, we use statistical learning techniques in a supervised strategy (prior  
356 classification of the samples) to build a model based on metabolomics data for  
357 identification of samples to species. Specifically, we used Random Forest to assess the  
358 overall accuracy of chemocoding to distinguish between the 82 species of *Inga* that we  
359 sampled at the five study sites (five saplings per species per site). Random Forest works  
360 by creating many classification trees each trained using random bootstrapped samples  
361 from the original metabolomics data. A consensus classification is then chosen based on  
362 the majority vote from all trees (Breiman *et al.*, 1984). Out of the five samples per  
363 species we iteratively dropped one sample to train the model with four samples and test  
364 the predictive accuracy with the fifth sample, such that each sample was used once for  
365 testing and four times for training.

366 The resulting Random Forest model accurately classifies 94% of the 410  
367 individuals representing the 82 *Inga* species (each with representatives from a single site,  
368 Table S2). Twenty-six of these species occurred at two to four of the study sites, so we  
369 also examined the model's classification accuracy when regional variation across sites

370 was included. In this case, our analyses correctly classified samples to species 96% of the  
371 time (Table S3). In addition, 90% of the time the classification model identified the  
372 correct species and site, indicating that there were regional differences in secondary  
373 metabolites within a species (Table S3). Overall, these results demonstrate that even in  
374 the face of cross-site intraspecific variation, there is sufficiently high interspecific  
375 variation to allow unknown samples to be efficiently classified into units that correspond  
376 to species as defined by morphology and DNA (Table S2).

377

## 378 **Discussion**

379 *The need for new tools.* The urgent need to catalogue, manage and understand the  
380 ecology and evolution of plant diversity has led to the development of innovative tools to  
381 improve the discrimination and identification of plant species. Traditional morphological-  
382 and molecular-based taxonomic identification methods have proved problematic for  
383 species-rich regions and highly species-rich genera (Dick & Webb, 2012; Seberg &  
384 Petersen, 2009). And it is precisely these diverse situations where accurate species  
385 identifications are most crucial. Alternative new technologies for plant identification in  
386 tropical trees include the use of near-infrared (NIR) leaf spectroscopy. However, its  
387 potential as a taxonomic tool has not been assessed across broad geographic scales for  
388 widespread species (Dugarte *et al.*, 2013; Lang *et al.*, 2013; Baker *et al.*, 2017).

389 Here, we add chemocoding to the tool box, and show that small, defense-related  
390 chemical markers characterized via untargeted metabolomics have great potential for  
391 species identification and in providing additional evidence for species delimitation and  
392 taxonomic discovery. Chemocoding was a robust method for species identification at a  
393 single site and across broad geographic scales, even for *Inga*, where levels of  
394 interspecific morphological variation are low and DNA barcoding is ineffective (Kress *et*  
395 *al.*, 2009; Gonzalez *et al.*, 2009; Dexter *et al.*, 2010).

396 *Accuracy of identifications.* An effective species identification method must be  
397 diagnostic of species (recognizing all populations rather than only sub-specific units or  
398 populations), and involve traits that are always present (rather than inducible or otherwise  
399 phenotypically plastic). Our analysis of *Inga* shows that entities identified as discrete  
400 morphospecies or phylogenetic taxa can indeed show constitutive differences in chemical

401 defenses that result in distinct chemocodes. For example, our results show that  
402 chemocoding correctly identified species at a single site and, most importantly, across  
403 their ranges (Figs. 2-5, Figs. S2-S9, Table S3). Even with 410 individuals from 82  
404 species of *Inga*, with a given species occurring at multiple sites, the Random Forest  
405 Analysis based on chemocodes accurately classified 96% of the individuals (Table S2). It  
406 also appears to be a robust method for identifying widespread species where intraspecific  
407 geographical variation might be problematic. For example, chemocoding of the widely  
408 scattered populations of *I. auristellae* resolves them as a single group, which is also  
409 resolved as monophyletic in our DNA sequence-based phylogeny (Fig. 5). Nevertheless,  
410 it is important to consider that depending on levels of migration, polymorphism, and  
411 selection, other outcomes are possible. In particular, some species that are widespread  
412 may show significant divergence in chemistry (e.g. *I. alata* (Fig. S8), and *I. marginata*  
413 (Fig. S9)). Hence, we caution that the efficacy of chemocode-based identification should  
414 be explored in each candidate taxon. However, we conclude that given the abundance and  
415 diversity of *Inga* species in neotropical forests, and the difficulty of identifying them  
416 using morphological characters (particularly in sterile material), chemocodes provide a  
417 valuable taxonomic tool.

418 Chemocoding is unlikely to identify species reliably where major components of  
419 plant chemistry show phenotypic plasticity. Rather than being constitutive, chemistry  
420 could be age-dependent (ontogeny or tissue age) or inducible by herbivores, pathogens,  
421 light, etc. Since this could generate significant within-site variation, we recognize that  
422 chemocoding may not work for species where important chemical markers vary.  
423 However, studies with several species-rich genera in the tropics have found that inter-  
424 specific differences in the defensive metabolome are large relative to intra-specific  
425 variation, even considering factors that are recognized as generators of plastic variation  
426 such as leaf ontogeny (expanding vs. mature leaves; Sedio *et al.*, 2017, Wiggins *et al.*,  
427 2016), light environment (sun vs. shade; Sinimbu *et al.*, 2012; Bixenmann *et al.*, 2016;  
428 Sedio *et al.*, 2017), season (dry vs. wet; Sedio *et al.*, 2017), and induction by herbivory  
429 (Bixenmann *et al.*, 2016). Even though plasticity may be an issue in some taxa, it does  
430 not rule out useful application of chemocoding, but highlights the need to separate  
431 diagnostic from phenotypically plastic characters.

432           *Practicality.* An identification method also needs to be practical. In other words, it  
433 needs to be accurate, rapid and inexpensive, as is the case with DNA barcoding.  
434 However, in groups where barcoding using standard markers cannot discriminate among  
435 species, sequencing of many genes may be necessary, which can be time-consuming and  
436 expensive. In the case of *Inga*, obtaining the resolved phylogeny took several years and  
437 hundreds of thousands of dollars (Nicholls *et al.*, 2015). Furthermore, only a few  
438 individuals of each species were sequenced. However, as per-base pair prices for  
439 sequencing are dropping in next-generation sequencing approaches, discriminating  
440 among species based on DNA sequences from many hundreds of loci may become more  
441 feasible.

442           For widespread use of chemocoding, we envision the creation of a public library  
443 of reference “chemocodes”, analogous to iBOL (ibol.org). In principal, chemocodes  
444 could be similar to barcodes in that they employ a limited set of compounds that are both  
445 constitutively produced and diagnostic of species. Instead, our method relies on the entire  
446 chemical fingerprint (typically a suite of more than 100 compounds) to identify species.  
447 Our choice is based on the variation observed single species within and across sites (e.g.  
448 *I. auristellae*, Fig. 5), taking into account variation caused by ontogeny.

449           We have successfully tested chemocoding on taxa from other groups, such as  
450 species from the families Euphorbiaceae, Malvaceae, Moraceae, Rubiaceae, Violaceae,  
451 among others (data not shown), suggesting that our approach works with groups other  
452 than *Inga*. Our data are publicly available (see Material and Methods) allowing others to  
453 attempt to develop diagnostic compounds for species identification. Nevertheless, the  
454 most challenging issue to address before chemocoding can be widely used, will be the  
455 application of our approach across different laboratories. For this, one must ensure that  
456 the same metabolic traits are used in all laboratories. In contrast to DNA barcoding,  
457 which uses standardized markers across taxa, each species is scored using different traits.  
458 At present, we do not know the extent to which chemical fingerprints will differ,  
459 depending for example on the exact column used for liquid chromatography or the exact  
460 model of mass spectrometer used. To address this issue such that chemocoding can be  
461 applied generally may require a more rigorous approach, in particular, the application of  
462 tandem MS or MS/MS (instead of simplifying the analysis as suggested below). Another

463 issue is that, while we used compounds with intermediate polarity for chemocoding of  
464 *Inga*, non-polar compounds such as terpenes or highly polar compounds such as non-  
465 protein amino acids may work best with other clades.

466 Because ecological studies based on long-term monitoring plots require the  
467 accurate identification of thousands of juvenile and adult trees, we consider here whether  
468 chemocoding could be applied to large numbers of samples. Currently, we have used  
469 chemocoding rather than DNA-based analyses for classifying over one thousand samples  
470 of *Inga* and have found chemocoding to be effective and convenient.

471 In terms of scaling up, one consideration is that these analyses can be carried out  
472 more expediently. The methods used in the present study range from easily accomplished  
473 to more complex methods. For example, sample preparation (e.g. drying) in the field and  
474 sample extraction in the laboratory are both straightforward. In addition, chemical  
475 analysis is largely free of contamination issues, which are a serious concern in DNA  
476 barcoding analyses, where specificity of primers may be low (Hollingsworth *et al.*, 2011).  
477 Other components could be simplified to streamline the analysis. These include collecting  
478 mature instead of expanding leaves. Most often, rainforest plants do not have expanding  
479 leaves, restricting chemocoding to a minority of saplings at any given point in time and  
480 reducing the utility of chemocoding. We found that mature leaves have most, but not all,  
481 of the chemical signals found in expanding leaves (Lokvam *et al.*, 2007; Wiggins *et al.*,  
482 2016), so the use of mature leaves should be feasible. Additionally, we used UPLC with a  
483 150mm column, followed by detection with a high-resolution time-of-flight mass  
484 spectrometer, an expensive analysis. To simplify this, we recommend using a shorter,  
485 50mm column, saving solvent and instrument time, followed by detection with a  
486 quadrupole mass spectrometer. Single quadrupole, triple quadrupole or ion trap detectors  
487 are much less expensive than a time-of-flight. While these provide lower mass resolution,  
488 both negative and positive mode data can be obtained in a single run, something that  
489 generally is not possible with a time-of-flight spectrometer. Based on our extensive work  
490 on *Inga*, only some of which is presented here, most pairs of similar species should be  
491 distinguishable using the proposed simplifications. There are some cases where two  
492 species have similar chemistry and morphology and the simpler chemical analyses may  
493 lump these into a single chemotype. But our experience suggests that these would show

494 high “within-chemotype” variation, indicating the need for more sophisticated chemical  
495 methods that can effectively answer the question at hand. In our hands, chemocoding  
496 gave clear results and was practical in terms of time and cost. Using the simplifications  
497 suggested above would decrease cost and time. We estimate that manual extraction and  
498 automated chromatography and data analysis could take 20 min and \$20.00 per sample at  
499 the time of writing. In summary, because chemocoding may work in circumstances where  
500 barcoding does not, we propose that it presents a novel and practical approach for  
501 surveying large number of individuals, possibly thousands of samples.

502

### 503 **Conclusions**

504 Our aspiration is not to claim that chemocoding will replace DNA barcoding or  
505 morphology-based identification methods, nor that chemocoding can determine species  
506 boundaries. Instead we suggest that it is a tool that provides valuable additional source of  
507 data to facilitate identification of plant species, especially for groups where traditional  
508 methods fall short. Chemocoding may be especially valuable in identifying species in  
509 recent radiations where morphological distinctions between species are slight and  
510 standard barcode markers do not provide sufficient resolution. However, it could also be  
511 used more broadly for species identification in cases where hundreds or thousands of  
512 samples need analysis, with the added benefit of providing information on defensive  
513 metabolites. In general, it can help by standardizing species names across multiple sites,  
514 and even in pin-pointing entities that may be species new to science. Such is the case for  
515 our third example (case study 3: *I. obidensis* and *I.cf brachystachys* in Brazil, French  
516 Guiana and Peru), where chemocoding and DNA suggest that the samples collected in  
517 Peru might represent a new species. Our approach is especially amenable for field  
518 biologists who work in networks of forest inventory plots since it can consistently  
519 distinguish amongst multiple species across geographic space (Figs. 4, 5; Tables S2, S3),  
520 and hence, can help in taxonomic integration across plots.

521       Experience with some taxa may show that chemocodes could help to distinguish  
522 groups of individuals that show similar plastic responses to shared abiotic environments  
523 or natural enemies (regardless of whether these correspond to species, e.g. Fig. S5). If so,  
524 chemocoding could be a very valuable way of dividing individual plants into ecologically

525 significant sets and improving our understanding of the plant-herbivore adaptive  
526 landscape. Given the key role of plant chemistry in many aspects of plant-herbivore-  
527 enemy interactions, it may be tremendously valuable to see plant communities in terms of  
528 chemotypes of hosts experienced by herbivores (Endara *et al.*, 2017).

529 The use of metabolites as tools in systematics has a long history and many  
530 antecedents (Gibbs, 1974; Smith, 1976; Harborne & Turner, 1984). While this has  
531 traditionally been used to investigate evolutionary relationships, based on the presence,  
532 absence and distinctive structures of specific classes of secondary metabolites in different  
533 groups at all taxonomic levels (Singh, 2016), chemocoding differs in that it is designed to  
534 *quantitatively* discriminate species across samples. This uses an unsupervised statistical  
535 approach that will likely yield consistent results independently of a priori ideas of species  
536 classifications. We see great potential in chemocoding to assist in the inventory of  
537 species-rich forests and potentially in the discovery of new species.

538

### 539 **Acknowledgments**

540 We thank the Ministry of Environment of Ecuador and the Ministry of Agriculture of  
541 Peru for granting the research and exportation permits. Valuable field assistance was  
542 provided by Zachary Benavidez, Allison Thompson, Yamara Serrano and Mayra  
543 Ninazunta. This work was supported by grants from the National Science Foundation  
544 (DEB-0640630 and DIMENSIONS of Biodiversity DEB-1135733), and Nouragues  
545 Travel Grants Program, CNRS, France to T.A.K. and P.D.C., and the Secretaría Nacional  
546 de Educación Superior, Ciencia, Tecnología e Innovación del Ecuador (SENESCYT) to  
547 M.J.E.

548

### 549 **Author contributions**

550 M.J.E., P.D.C., N.L.W., D.L.F and T.A.K. designed and conducted the research. M.J.E.,  
551 N.L.W. and D.L.F. designed and performed the data analysis. G.C.Y. contributed to the  
552 metabolomic analysis. J.A.N., R.T.P., K.G.D., C.A.K. and G.N.S. contributed the next-  
553 generation DNA sequence data. M.J.E., P.D.C., N.L.W., D.L.F., J.A.N., R.T.P., K.G.D.,  
554 C.A.K., G.N.S. and T.A.K. wrote the manuscript.

555

556 **References**

- 557 Baker TR, Pennington RT, Dexter KG, Fine PVA, Fortune-Hopkins H, Honorio EN,  
558 Huamatunpa-Chuquimaco I, Klitgård BB, Lewis GP, de Lima HC, *et al.* 2017.  
559 Maximizing synergy among tropical plant systematists, ecologists, and evolutionary  
560 biologists. *Trends in Ecology & Evolution* **32**: 258-267.
- 561
- 562 Bartel J, Krumsiek J, Theis FJ. 2013. Statistical methods for the analysis of high-  
563 throughput metabolomics data. *Computational and Structural Biotechnology Journal* **4**:  
564 e201301009.
- 565
- 566 Benton P, Want EJ, Ebbels TMD. 2010. Correction of mass calibration gaps in liquid  
567 chromatography–mass spectrometry metabolomics data. *Bioinformatics* **26**:2488–  
568 2489.
- 569
- 570 Bixenmann RJ, Coley PD, Weinhold A, Kursar TA. 2016. Higher herbivore pressure  
571 favors constitutive over induced defense. *Ecology and Evolution* **6**: 6037-6049.
- 572
- 573 Breiman L, Friedman JH, Olshen RA, Stone CJ. 1984. *Classification and Regression*  
574 *Trees*. Monterey, CA, USA: Wadsworth.
- 575
- 576 Breiman L. 2001. Random forests. *Machine Learning* **45**: 5–32.
- 577
- 578 Chen S, Yao H, Han J, Liu C, Song J, Shi L, Zhu Y, Ma X, Gao T, Pang X, *et al.* 2010.  
579 Validation of the ITS2 region as a novel DNA barcode for identifying medicinal plant  
580 species. *PLoS ONE* **5**: e8613.
- 581
- 582 China Plant iBOL Group, Li DZ, Gao LM, Li HT, Wang H, Ge XJ, Liu JQ, Chen ZD,  
583 Zhou SL, Chen SL, *et al.* 2011. Comparative analysis of a large dataset indicates that  
584 internal transcribed spacer (ITS) should be incorporated into the core barcode for seed  
585 plants. *Proceedings of the National Academy of Sciences USA* **108**: 19641–19646.
- 586

587 Costion C, Ford A, Cross H, Crayn D, Harrington M, Lowe A. 2011. Plant DNA  
588 barcodes can accurately estimate species richness in poorly known floras. *PLoS ONE* **6**:  
589 e26841.  
590

591 Cristescu ME. 2014. From barcoding single individuals to metabarcoding biological  
592 communities: towards an integrative approach to the study of global biodiversity. *Trends*  
593 *in Ecology & Evolution* **29**: 566-571.  
594

595 Cutler DR, Edwards KH, Beard AC, Kyle TH, Jacob G, Joshua JL. 2007. Random forests  
596 for classification in ecology. *Ecology* **88**: 2783–2792.  
597

598 Cutler A, Cutler DR, Stevens JR. 2009. Tree-based methods. In: Li X, Xu R, eds. *High*  
599 *dimensional data analysis in cancer research*. New York, USA: Springer-Verlag, 83-101.  
600

601 de Carvalho MR, Bockmann FA, Amorim DS, de Vivo M, de Toledo-Piza M, Menezes  
602 NA, de Figueiredo JL, Castro R, Gill AC, McEachran JD. 2005. Revisiting the taxonomic  
603 impediment. *Science* **307**: 353.  
604

605 Dexter KG, Pennington TR, Cunningham CW. 2010. Using DNA to assess errors in  
606 tropical tree identifications: How often are ecologists wrong and when does it matter?  
607 *Ecological Monographs* **80**: 267-286.  
608

609 Dick CW, Webb CO. 2012. Plant DNA barcodes, taxonomic management and species  
610 discovery in tropical forests. In: Kress WJ, Erickson DL, eds. *DNA barcodes: Methods*  
611 *and protocols*. Totowa, NJ: Humana Press, 379-383.  
612

613 Dugarte FM, Higuchi N, Almeida A, Vicentina A. 2013. Species spectral signature:  
614 Discriminating closely related species in the Amazon with Near-Infrared Leaf-  
615 Spectroscopy. *Forest Ecology and Management* **291**: 240-248.  
616

617 Endara MJ, Coley PD, Ghabash G, Nicholls JA, Dexter KG, Donoso DA, Stone GN,  
618 Pennington RT, Kursar TA. 2017. Coevolutionary arms race versus host defense chase in  
619 a tropical-herbivore plant system. *Proceedings of the National Academy of Sciences USA*,  
620 **114**: E7499–E7505.

621

622 Embrechts MJ, Gatti CJ, Linton J, Roysam B. 2013. Hierarchical Clustering for large  
623 data sets. In: Georgieva P, Mihaylova L, Jain LC, eds. *Advances in intelligent signal*  
624 *processing and data mining: Theory and applications*. New York, USA: Springer, 197-  
625 233.

626

627 Gemeinholzer B, Oberprieler C, Bachmann K. 2006. Using GenBank data for plant  
628 identification: possibilities and limitations using the ITS 1 of Asteraceae species  
629 belonging to the tribes Lactuceae and Anthemideae. *Taxon* **55**: 173-187.

630

631 Gibbs RD. 1974. *Chemotaxonomy of flowering plants*. Montreal and London: McGill's  
632 Queen's University Press.

633

634 Gonzalez MA, Baraloto C, Engel J, Mori SA, Pétronelli P, Riéra B, Roger A, Thébaud C,  
635 Chave J. 2009. Identification of Amazonian trees with DNA barcodes. *PLoS ONE* **4**:  
636 e7483.

637

638 Grimaldi C, Riéra B. 2001. Geography and climate. In: Bongers F, Charles-Dominique,  
639 P, Forget PM, Théry, M, eds. *Nouragues: Dynamics and plant-animal interactions in a*  
640 *Neotropical rain forest*. New York, USA: Springer-Science + Business Media.

641

642 Harborne JB, Turner BL. 1984. *Plant chemosystematics*. London, UK: Academic Press.

643

644 Hebert PDN, Cywinska A, Ball SL, deWaard JR. 2003. Biological identifications through  
645 DNA barcodes. *Proceedings of the Royal Society B: Biological Sciences* **270**: 313-321.

646

647 Hollingsworth PM, Forrest LL, Spouge JL, Hajibabaei M, Ratnasingham S, van der Bank  
648 M, Chase MW, Cowan RS, Erickson DL, Fazekas AJ, *et al.* 2009. A DNA barcode for  
649 land plants. *Proceedings of the National Academy of Sciences USA* **106**: 12794–12797.  
650

651 Hollingsworth PM, Graham SW, Little DP. 2011. Choosing and using a plant DNA  
652 barcode. *PLoS ONE* **6**: e19254.  
653

654 Kress WJ, Erickson DL, Jones FA, Swenson NG, Perez R, Sanjur O, Bermingham E.  
655 2009. Plant DNA barcodes and a community phylogeny of a tropical forest dynamics plot  
656 in Panama. *Proceedings of the National Academy of Sciences USA* **106**: 18621-18626.  
657

658 Kress WJ, Erickson DL, Swenson NG, Thompson J, Uriarte M, Zimmerman JK. 2010.  
659 Advances in the use of DNA barcodes to build a community phylogeny for tropical trees  
660 in a Puerto Rican forest dynamics plot. *PLoS ONE* **5**: e15409.  
661

662 Kress WJ, García-Robledo C, Uriarte M, Erickson DL. 2015. DNA barcodes for ecology,  
663 evolution, and conservation. *Trends in Ecology & Evolution* **30**: 25-35.  
664

665 Kuhl C, Tautenhahn R, Bottcher C, Larson TR, Neumann S. 2012. CAMERA: an  
666 integrated strategy for compound spectra extraction and annotation of liquid  
667 chromatography/mass spectrometry data sets. *Analytical Chemistry* **84**: 283-289.  
668

669 Kursar TA, Dexter KG, Lokvam J, Pennington RT, Richardson JE, Weber MG,  
670 Murakami ET, Drake C, McGregor R, Coley PD. 2009. The evolution of antiherbivore  
671 defenses and their contribution to species coexistence in the tropical tree genus *Inga*.  
672 *Proceedings of the National Academy of Science USA* **106**:18073–18078.  
673

674 Lang C, Costa FRC, Camargo JLC, Durgante M, Vicentini A. 2015. Near infrared  
675 spectroscopy facilitates rapid identification of both young and mature Amazonian tree  
676 species. *PLoS ONE* **10**: e0134521.  
677

678 Lawrence RL, Shana DW, Roger LS. 2006. Mapping invasive plants using hyperspectral  
679 imagery and Breiman Cutler Classifications (randomForest). *Remote Sensing of*  
680 *Environment* **100**: 356–362.

681

682 Leigh EG. 1999. *Tropical forest ecology: a view from Barro Colorado Island*. New York,  
683 USA: Oxford University Press.

684

685 Liaw A, Wiener M. 2002. Classification and regression by randomForest. *R News* 2: 18-  
686 22.

687

688 Liu J, Yan HF, Newmaster SG, Pei N, Ragupathy S, Ge XJ. 2015. The use of DNA  
689 barcoding as a tool for the conservation biogeography of subtropical forests in China.  
690 *Diversity of Distributions* **21**: 188-199.

691

692 Lokvam J, Clausen TP, Grapov D, Kursar TA. 2007. Galloyl depsides of tyrosine from  
693 young leaves of *Inga laurina*. *Journal of Natural Products* **70**: 134-136.

694

695 Mace GM. 2004. The role of taxonomy in species conservation. *Philosophical*  
696 *Transactions of the Royal Society of London Series B, Biological Sciences* **359**: 711–719.

697

698 Martucci MEP, De Vos RCH, Carollo CA, Gobbo-Neto L. 2014. Metabolomics as a  
699 potential chemotaxonomical tool: application in the genus *Vernonia* Schreb. *PLoS ONE* **9**:  
700 e93149.

701

702 Naciri Y, Linder HP. 2015. Species delimitation and relationships: the dance of the seven  
703 veils. *Taxon* **64**: 3–16.

704

705 Nicholls JA, Pennington RT, Koenen EJM, Hughes CE, Hearn J, Bunnefeld L, Dexter  
706 KG, Stone GN, Kidner CA. 2015. Using targeted enrichment of nuclear genes to increase  
707 phylogenetic resolution in the neotropical rainforest genus *Inga* (Leguminosae:  
708 Mimosoideae). *Frontiers in Plant Science* **6**: 1-20.

709

710 Pennington TD. 1997. *The Genus Inga: Botany*. UK: The Royal Botanic Gardens, Kew.

711

712 Pennington TD, Revelo N. 1997. *El género Inga en el Ecuador*. UK: The Royal Botanic  
713 Gardens, Kew.

714

715 Pennington RT, Lavin M. 2015. The contrasting nature of woody plant species in  
716 different neotropical forest biomes reflects differences in ecological stability. *New*  
717 *Phytologist* **210**: 25-37.

718

719 Pitman N. 2007. An overview of the Los Amigos watershed, Madre de Dios, southeastern  
720 Peru. October 2007 version of an unpublished report available from the author at  
721 [npitman@amazonconservation.org](mailto:npitman@amazonconservation.org)

722

723 Radtke MG, Da Fonseca CRV, Williamson GB. 2007. The old and young Amazon: Dung  
724 beetle biomass, abundance and species diversity. *Biotropica* **39**: 725–730.

725

726 Razafimandimbison SG, Kellogg EA, Bremer B. 2004. Recent origin and phylogenetic  
727 utility of divergent ITS putative pseudogenes: a case study from Naucleaeae (Rubiaceae).  
728 *Systematic Biology* **53**: 177-192.

729

730 Reeb PD, Bramardi SJ, Steibel JP. 2015. Assessing dissimilarity measures for sample-  
731 based hierarchical clustering of RNA sequencing data using plasmid datasets. *PloS*  
732 *ONE* **10**: e0132310.

733

734 Richardson JE, Pennington RT, Pennington TD, Hollingsworth PM. 2001. Rapid  
735 diversification of a species-rich genus of neotropical rain forest trees. *Science* **293**: 2242–  
736 2245.

737

738 Seberg O, Petersen G. 2009. How many loci does it take to DNA barcode a *Crocus*? *PloS*  
739 *ONE* **4**: e4598.

740

741 Sedio BE, Rojas Echeverri JC, Boya CA, Wright JS. 2017. Sources of variation in foliar  
742 secondary chemistry in a tropical forest tree community. *Ecology* **98**: 616-623.

743

744 Singh R. 2016. Chemotaxonomy: a tool for plant classification. *Journal of Medicinal*  
745 *Plant Studies* **4**: 90-93.

746

747 Sinimbu G, Coley PD, Lemes MR, Lokvam J, Kursar TA. 2012. Do the antiherbivore  
748 traits of developing leaves in the Neotropical tree *Inga paraensis* (Fabaceae) vary with  
749 light availability? *Oecologia* **170**: 669-676.

750

751 Smith PM. 1976. *The chemotaxonomy of plants*. London, UK: Edward Arnold.

752

753 Smith CA, Want EJ, O'Maille G, Abagyan R, Siuzdak G. 2006. XCMS: processing mass  
754 spectrometry data for metabolite profiling using nonlinear peak alignment, matching, and  
755 identification. *Analytical Chemistry* **78**: 779–787.

756

757 Suzuki R, Shimodaira H. 2014. pvclust: hierarchical clustering with p-values via  
758 multiscale bootstrap resampling. *Bioinformatics* **22**:1540–1542.

759

760 Tautenhahn R, Böttcher C, Neumann S. 2008. Highly sensitive feature detection for high  
761 resolution LC/MS. *BMC Bioinformatics* **9**:504.

762

763 Tautenhahn R, Patti GJ, Kalisiak E, Miyamoto T, Schmidt M, Lo FY, McBee J, Baliga  
764 NS, Siuzdak G. 2011. MetaXCMS: second-order analysis of untargeted metabolomics  
765 data. *Analytical Chemistry* **83**: 696-700.

766

767 ter Steege H, Pitman NCA, Sabatier D, Baraloto C, Salomão RP, Guevara JE, Phillips  
768 OL, Castilho CV, Magnusson WE, Moline JF, *et al.* 2013. Hyperdominance in the  
769 Amazonian tree flora. *Science* **342**: 1243092.

770

771 Valencia R, Balslev H, Miño GP. 1994. High tree  $\alpha$  diversity in Amazonian  
772 Ecuador. *Biodiversity and Conservation* **3**: 21–28.  
773  
774 Valencia R, Foster RB, Villa G, Condit R, Svenning J-C, Hernández C, Romoleroux K,  
775 Losos E, Magård E, Balslev H. 2004. Tree species distributions and local habitat  
776 variation in the Amazon: large forest plot in eastern Ecuador. *Journal of Ecology* **92**:  
777 214-229.  
778  
779 Valentini A, Pompanon F, Taberlet P. 2009. DNA barcoding for ecologists. *Trends in*  
780 *Ecology & Evolution* **24**: 110-117.  
781  
782 Vernooij R, Haribabu E, Muller MR, Vogel JH, Hebert PDN, Schindel DE, Shimura J,  
783 Singer GAC. 2010. Barcoding life to conserve biological diversity: beyond the taxonomic  
784 imperative. *PLoS Biology* **8**: e1000417.  
785  
786 Wiggins NL, Forrister DL, Endara MJ, Coley PD, Kursar TA. 2016. Quantitative and  
787 qualitative shifts in defensive metabolites define chemical defense investment during leaf  
788 development in *Inga*, a genus of tropical trees. *Ecology and Evolution* **6**: 478-492.

## 790 **Figure legends**

791  
792 **Figure 1.** Study sites: 1) Barro Colorado, Panamá, 2) Nouragues, French Guiana, 3)  
793 Tiputini, Ecuador, 4) KM41 near Manaus, Brazil and, 5) Los Amigos, Perú.  
794

795 **Figure 2.** Case study 1: Two morphologically confusing or cryptic species that are  
796 present at the same site: *I. alata* and *I. pezizifera* in French Guiana. a) Total ion  
797 chromatograms showing relative intensities of peaks from the LC-QToF-MS in negative  
798 mode, b) Hierarchical cluster dendrograms based on relative abundances of UPLC-MS  
799 metabolites. The numbers in red above each branch point are the Approximately  
800 Unbiased confidence levels, these indicate the probability that the samples below that  
801 point are a cluster. Clusters with values of 95 signify  $P=0.05$ , indicating that these

802 clusters are strongly supported by the data, c) Clade containing *I. alata* and *I. pezizifera*  
803 adapted from a resolved phylogeny based on next-generation DNA sequence data  
804 (Nicholls *et al.*, 2015 & unpubl. data). Numbers in black represent bootstrap support  
805 values. Values greater than 95 indicate that the clade is strongly supported by the data.

806

807 **Figure 3.** Case study 2: A single species that shows morphological variation within a  
808 site: *I. acreana* T28 and *I. acreana* T56 in Ecuador. a) Total ion chromatograms showing  
809 relative intensities of peaks from the LC-QToF-MS in negative mode, b) Hierarchical  
810 clustering based on relative abundances of UPLC-MS metabolites. The numbers in red  
811 above each branch point are the Approximately Unbiased confidence levels, these  
812 indicate the probability that the samples below that point are a cluster. Clusters with  
813 values of 95 signify  $P=0.05$ , indicating that these clusters are strongly supported by the  
814 data, c) Clade containing *I. acreana* T28 and *I. acreana* T56 adapted from a resolved  
815 phylogeny based on next generation DNA sequence data (Nicholls *et al.*, 2015 & unpubl.  
816 data).

817

818 **Figure 4.** Case study 3: Distinguishing among two morphologically similar species  
819 across three sites: *I. obidensis* from Brazil and French Guiana and *I. brachystachys* from  
820 Peru. a) Total ion chromatograms showing relative intensities of peaks from the LC-  
821 QToF-MS in negative mode, b) Hierarchical clustering based on relative abundances of  
822 UPLC-MS metabolites. The numbers in red above each branch point are the  
823 Approximately Unbiased confidence levels, these indicate the probability that the  
824 samples below that point are a cluster. Clusters with values of 95 signify  $P=0.05$ ,  
825 indicating that these clusters are strongly supported by the data, c) Clade containing *I.*  
826 *obidensis* and *I. cf brachystachys* adapted from a resolved phylogeny based on next-  
827 generation DNA sequence data (Nicholls *et al.*, 2015 & unpubl. data).

828

829 **Figure 5.** Case study 4: Identification of a widespread species across its range: *I.*  
830 *auristellae*. a) Total ion chromatograms showing relative intensities of peaks from the  
831 LC-QToF-MS in negative mode, b) Hierarchical clustering based on relative abundances  
832 of UPLC-MS metabolites. The numbers in red above each branch point are the

833 Approximately Unbiased confidence levels, these indicate the probability that the  
834 samples below that point are a cluster. Clusters with values of 95 signify  $P=0.05$ ,  
835 indicating that these clusters are strongly supported by the data, c) Clade containing *I.*  
836 *auristellae* adapted from a resolved phylogeny based on next-generation DNA sequence  
837 data (Nicholls *et al.*, 2015 & unpubl. data).

838

### 839 **Supporting Information**

840

841 **Fig. S1.** Photographs of the *Inga* species for each study case.

842

843 **Fig. S2.** Two morphologically confusing or cryptic species that are present at the same  
844 site: *Inga coruscans* and *I. laurina* in Peru.

845

846 **Fig. S3.** Two morphologically confusing or cryptic species that are present at the same  
847 site: *Inga umbellifera* and *I. microcoma* in Ecuador.

848

849 **Fig. S4.** Two morphologically confusing or cryptic species that are present at the same  
850 site: *Inga chartacea* and *I. sapindoides* in Ecuador.

851

852 **Fig. S5.** Morphological and chemical variation within a site: *Inga leiocalycina* T65 and *I.*  
853 *leiocalycina* T86 in Ecuador.

854

855 **Fig. S6.** Identification of a widespread species across its range: *Inga pezizifera* in  
856 Panama, French Guiana and Brazil.

857

858 **Fig. S7.** Identification of a widespread species across its range: *Inga alba* in French  
859 Guiana, Brazil and Peru.

860

861 **Fig. S8.** Identification of a widespread species across its range: *Inga alata* in French  
862 Guiana, Ecuador and Peru.

863

864 **Fig. S9.** Identification of a widespread species across its range: *Inga marginata* in  
865 Panama, French Guiana, Ecuador and Peru.

866

867 **Table S1.** Mass difference rules used for peak annotation in negative mode with  
868 CAMERA.

869

870 **Table S2.** Species accuracy classification from Random Forest Analyses for *Inga* species  
871 with representatives from a single site.

872

873 **Table S3.** Species accuracy classification from Random Forest Analyses for *Inga* species  
874 that occurred at two to four of the study sites.

875

876 **Methods S1.** R code for LC/MS raw data pre-processing in XCMS

877

878 **Methods S2.** R code for Random Forest Analysis

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894



**Table 1.** Study species

<b>Figures</b>	<b>Case study</b>	<b>Species and sites</b>
Figure 2	Two morphologically confusing species within a site	<i>Inga alata</i> Benoist, and <i>Inga pezizifera</i> Benth, French Guiana
Figure 3	Morphological variation within one species at a site	<i>Inga acreana</i> <sup>1</sup> Harms, Ecuador
Figure 4	Two morphologically similar species across sites	<i>Inga</i> cf. <i>brachystachys</i> Ducke, and <i>Inga obidensis</i> Ducke, French Guiana, Brazil and Peru
Figure 5	Identification of a widespread species across its range	<i>Inga auristellae</i> Harms, Ecuador, Brazil and Peru
Figure S1	Two morphologically confusing species within a site	<i>Inga coruscans</i> Humb. & Bonpl. ex Willd., and <i>Inga laurina</i> (SW.) Willd., Peru
Figure S2	Two morphologically confusing species within a site	<i>Inga microcoma</i> Harms and <i>Inga umbellifera</i> <sup>2</sup> (Vahl) Steud., Ecuador
Figure S3	Two morphologically confusing species within a site	<i>Inga chartacea</i> Poepp. and <i>Inga sapindoides</i> Willd., Ecuador
Figure S4	Morphological and chemical variation within on species at a site	<i>Inga leiocalycina</i> <sup>3</sup> Benth., Ecuador

**Table 1.** Continued...

<b>Figure</b>	<b>Case study</b>	<b>Species and site</b>
Figure S5	Identification of a widespread species across its range	<i>Inga alata</i> Benoist, French Guiana, Ecuador and Peru
Figure S6	Identification of a widespread species across its range	<i>Inga pezizifera</i> Benth, Panama, French Guiana and Brazil
Figure S7	Identification of a widespread species across its range	<i>Inga alba</i> (SW.) Willd., French Guiana, Brazil and Peru
Figure S8	Identification of a widespread species across its range	<i>Inga marginata</i> Willd. Panama, French Guiana, Ecuador and Peru

896 <sup>1</sup>*I. acreana* includes two morphotypes in Ecuador: T28 and T56.

897 <sup>2</sup>*I. umbellifera* includes two morphotypes in Ecuador: T50 and T73.

898 <sup>3</sup>*I. leiocalycina* includes two morphotypes in Ecuador: T65 and T86.

899