

1 **Estimating belowground plant abundance with DNA metabarcoding**

2 Silvia Matesanz<sup>1\*</sup>, David S. Pescador<sup>1</sup>, Beatriz Pías<sup>2</sup>, Ana M. Sánchez<sup>1</sup>, Julia Chacón-  
3 Labella<sup>3</sup>, Angela Illuminati<sup>1</sup>, Marcelino de la Cruz<sup>1</sup>, Jesús López-Angulo<sup>1</sup>, Adrián  
4 Escudero<sup>1</sup>

5 <sup>1</sup> Área de Biodiversidad y Conservación, Universidad Rey Juan Carlos. Tulipán s/n,  
6 28933, Móstoles, Spain

7 <sup>2</sup> Departamento de Biodiversidad, Ecología y Evolución. Universidad Complutense de  
8 Madrid. José Antonio Nováis 2, 28040, Madrid, Spain.

9 <sup>3</sup> Departamento de Medio Ambiente. Instituto Nacional de Investigación y Tecnología  
10 Agraria y Alimentaria (INIA) Ctra. de la Coruña Km 7.5, 28040 Madrid, Spain

11 \* Author for correspondence: [silvia.matesanzgarcia@gmail.com](mailto:silvia.matesanzgarcia@gmail.com)

12

13

14 **Running title:** Estimating root biomass with metabarcoding

15 **Abstract**

16 Most work on plant community ecology has been performed aboveground, neglecting  
17 the processes that occur in the soil. DNA metabarcoding, where multiple species are  
18 computationally identified in bulk samples, can help overcome the logistical limitations  
19 involved in sampling plant communities belowground. A major limitation of this  
20 methodology is, however, the quantification of species' abundances based on the  
21 percentage of sequences assigned to each taxon. Using root tissues of the five dominant  
22 species in a semiarid Mediterranean shrubland (*Bupleurum fruticosens*, *Helianthemum*  
23 *cinereum*, *Linum suffruticosum*, *Stipa pennata* and *Thymus vulgaris*), we built pairwise  
24 mixtures of relative abundance (20, 50 and 80% biomass), and implemented two  
25 methods (linear models fits and correction indices) to improve root biomass estimates.  
26 We validated both methods with multispecies mixtures that simulate field-collected  
27 samples. For all species, we found a positive and highly significant relationship between  
28 the percentage of sequences and biomass in the mixtures ( $R^2 = 0.44-0.66$ ), but the  
29 equations for each species (slope and intercept) differed among them, and two species  
30 were consistently over- and under-estimated. The correction indices greatly improved  
31 the estimates of biomass percentage for all five species in the multispecies mixtures,  
32 and reduced the overall error from 17% to 6%. Our results show that, through the use of  
33 post-sequencing quantification methods on mock communities, DNA metabarcoding  
34 can be effectively used to determine not only species' presence but also their relative  
35 abundance in field samples of root mixtures. Importantly, knowledge on these aspects  
36 will allow to study key, yet poorly understood, belowground processes.

37

38 **Keywords:** DNA metabarcoding, plant abundance, root biomass, sequence,  
39 Mediterranean shrubland, coexistence, mock communities, *rbcL* region

## 40 **1. INTRODUCTION**

41 A critical question in plant ecology is how communities are structured in space and  
42 time. In this still-unresolved debate, community ecologists attempt to disentangle the  
43 relative role of key stochastic and deterministic processes, such as niche differentiation,  
44 biotic interactions, and environmental filtering to determine plant species coexistence  
45 (Chase & Leibold 2003; Götzenberger et al. 2012; Gravel et al. 2006; HilleRisLambers  
46 et al. 2012; Vellend 2010). A major limitation is the fact that our understanding on the  
47 structure of plant diversity stems from data collected almost entirely aboveground.  
48 However, a large proportion of the community biomass can be located belowground  
49 (Hilbert & Canadell 1995; Poorter et al. 2012; Schenk & Jackson 2002), particularly in  
50 stressful habitats, and as such, both plant-soil and plant-plant interactions may have  
51 important implications for community-level processes (Bardgett et al. 2014; Bever et al.  
52 2010; Casper & Jackson 1997; Philippot et al. 2013; Wardle et al. 2004).

53         The main constraint to sampling plant communities belowground is that reliable  
54 species identification in natural conditions based solely on morphological root traits is  
55 extremely difficult or simply unfeasible in many cases (Silva & Rego 2003). In this  
56 context, the development of molecular methods such as DNA metabarcoding, spurred  
57 by the emergence of next-generation sequencing, has had a significant impact on  
58 biodiversity assessments (Schuster 2007; Taberlet et al. 2012). DNA metabarcoding  
59 involves the simultaneous identification of multiple species based on the amplification  
60 and sequencing of a common target DNA region from an environmental or community  
61 bulk sample (Deiner et al. 2017; Hollingsworth et al. 2009; Kress et al. 2005; Taberlet  
62 et al. 2012). For instance, for plant communities, DNA metabarcoding has been  
63 successfully used to recreate their current and past composition from soil-derived DNA  
64 (Fahner et al. 2016; Jørgensen et al. 2012; Porter et al. 2016; Yoccoz et al. 2012) or to

65 identify the floral composition of honey (Hawkins et al. 2015). Going belowground, a  
66 few studies have also assessed the richness and composition of temperate and tropical  
67 plant communities using root mixtures or individual root fragments (Hiiesalu et al.  
68 2012; Jones et al. 2011; Kesanakurti et al. 2011).

69         There is mounting evidence that DNA metabarcoding is a robust method to  
70 assess biodiversity. Indeed, some studies even found higher DNA-based diversity  
71 compared to traditional sampling methods (reviewed in Deiner et al. 2017). However,  
72 there is currently an intense debate on the use of read number to quantify DNA  
73 metabarcoding results, with some authors limiting its use to strictly detect occurrence,  
74 whilst others advocate a quantitative approach (see discussion on e.g. Bell et al. 2019;  
75 Deiner et al. 2017; Fonseca 2018; Porter & Hajibabaei 2018). Ideally, the percentage of  
76 sequences assigned to each taxa during DNA metabarcoding would closely reflect the  
77 species' abundance (biomass, number of individuals, etc.) in the bulk sample. Building  
78 on this simple assumption, several studies have attempted the direct use of the observed  
79 percentage of DNA sequences to estimate species' abundances in communities of  
80 microbes (Amend et al. 2010), stoneflies (Elbrecht & Leese 2015), fish and amphibians  
81 (Evans et al. 2016; Pont *et al.* 2018), zooplankton (Harvey et al. 2017) and fungi  
82 (Merges et al. 2018). However, many factors operating during DNA extraction,  
83 amplification and sequencing as well as the inherently compositional nature of the data  
84 can alter the correspondence between the percentage of reads retrieved and the species'  
85 abundance (Cristescu 2014; Deiner et al. 2017; Elbrecht & Leese 2015; Gloor,  
86 Macklaim, Pawlowsky-Glahn, & Egozcue, 2017; Pawlak et al. 2015; Polz & Cavanaugh  
87 1998; Porter & Hajibabaei 2018). Indeed, studies where such correspondence is lacking  
88 suggest that the use of uncorrected, observed percentages may render strongly biased  
89 estimates of abundance (see e.g. Bell et al. 2019; Deagle et al. 2013; Lim et al. 2016),

90 and effort is now being devoted to the improvement of quantification methods (Thomas  
91 et al. 2016 and references therein; McLaren et al. 2019; Nichols et al. 2018; Piñol et al.  
92 2015). In this context, the use of mock communities, i.e. a defined mixture of tissues  
93 with known species composition and relative abundance (biomass), can be a useful tool  
94 to improve biomass estimates in DNA metabarcoding studies (see e.g. Thomas *et al.*  
95 2016 for a comprehensive example using prey fish mixtures).

96         Quantification of species' biomass through DNA metabarcoding can be critical  
97 in the study of belowground community structure. Compared to other plant  
98 communities, Mediterranean shrublands are highly diverse, and up to 80% biomass can  
99 appear belowground (Hilbert & Canadell 1995). In these water-limited systems,  
100 belowground plant-plant interactions can be equally important, or even more, than those  
101 occurring aboveground (Casper & Jackson 1997). However, experimental evidence on  
102 their direction, strength and correspondence to the interactions occurring aboveground  
103 is scarce (but see Armas & Pugnaire 2011). Furthermore, because species' abundances  
104 are markedly heterogeneous and leptokurtic, with a few very abundant species and  
105 many rare ones (Chacón-Labela et al. 2017; Chacón-Labela et al. 2016; McGill et al.  
106 2007), presence-absence data fails to accurately reflect the structure of the plant  
107 community. Therefore, in order to gain insights on the mechanisms that determine plant  
108 community structure and to build a global coexistence theory, we should expand our  
109 focus belowground and compare these patterns to those aboveground. To do this, we  
110 need robust information not only on the presence of species in the soil but also on their  
111 relative abundance across space.

112         In this study, we built mock communities with varying composition and  
113 abundance of five selected species from Mediterranean shrublands, and used a DNA  
114 metabarcoding approach on these root mixtures, to move beyond species detection and

115 estimate species' relative biomass. We implement two post-sequencing quantification  
116 methods. First, we fit linear models to assess whether the percentage of reads (DNA  
117 sequences) can be used to robustly estimate percentage of root biomass, and second, we  
118 compute correction indices that control for potential biases and improve the relationship  
119 between sequences and biomass percentages (see Thomas et al. 2016). In addition, to  
120 determine the possibility to apply our results to field-collected samples, we validate  
121 both methods with multispecies realistic samples. To our knowledge, this is the first  
122 study aimed at the improvement of a quantitative DNA metabarcoding approach in  
123 plants.

124

## 125 **2. MATERIALS AND METHODS**

### 126 **2.1 Plant community and species selection**

127 The study plant community is a species-rich semiarid Mediterranean shrubland  
128 established in limestone and gypsum soil in the central Iberian Peninsula. Perennial  
129 cover ranges from 40 to 60%, and is mainly dominated by small chamaephytes and  
130 grasses. It is a highly diverse community, with around 50 perennial species found at the  
131 local scale (e.g.  $\approx 8000$  individuals from 48 species in  $60 \text{ m}^2$ ; Chacón-Labela *et al.*  
132 2016). The distribution of individuals across species is highly heterogeneous, with a few  
133 species accounting for a high proportion of the total number of individuals. Given the  
134 disproportionate influence of the most abundant species, we selected the five most  
135 common species in the community for our study (Fig. 1): *Thymus vulgaris* L.  
136 (Lamiaceae), *Helianthemum cinereum* (Cav.) Pers. (Cistaceae), *Linum suffruticosum* L.  
137 (Linaceae), *Bupleurum fruticosum* L. (Apiaceae), and *Stipa pennata* L. (Poaceae). The  
138 selected species have different phylogenetic origins, life forms, and can account for as  
139 much as 65% of the total number of individuals in the community (data not shown).

140

## 141 **2.2 Sampling material and creation of root mixture mock communities**

142 We collected root samples in the shrublands near Orusco de Tajuña (Madrid, Spain,  
143 40°17'17.5"N 3°12'19.4"W). For each selected species, we uprooted 5-10 adult  
144 individuals with unequivocal taxonomic identification. All individuals were collected  
145 within 24h, bagged separately, stored in a cooler and immediately transferred to the lab  
146 at Universidad Rey Juan Carlos. Upon arrival, their root system was thoroughly washed  
147 and separated from the soil, and roots from all individuals of the same species were  
148 pooled and maintained in cool water until sample preparation. We created mock  
149 communities (hereafter mixtures) based on mixtures of root biomass, varying both the  
150 species composition, richness and the percentage of biomass of each species in each  
151 sample. Note that the use of root mixtures (community DNA) rather than DNA  
152 extracted from soil samples (environmental DNA *sensu* Deiner et al. 2017) allowed to  
153 quantify biomass of actively growing plants, avoiding the presence of persistent DNA  
154 from long-dead individuals (Baird & Hajibabaei, 2012).

155 The communities were created by cutting small pieces of roots (removed of  
156 excess water by patting them with paper towel) and weighing them in a Mettler Toledo  
157 MX5microbalance (1 µg precision; Mettler Toledo, Columbus, OH, USA) the same day  
158 of collection in the field. All mock communities contained 100 mg of fresh root  
159 biomass, and were immediately frozen at -80°C for later DNA metabarcoding analyses.  
160 We created two different types of mock communities:

161 1) Pairwise mixtures, with two species present in different proportions (20:80,  
162 50:50, or 80:20; all pairwise combinations with three replicates per type of community,  
163 N = 90 samples; Fig. 1). The pairwise mixtures were used to: i) determine the match  
164 between the percentage of biomass and the percentage of DNA sequences (hereafter

165 reads) obtained via linear model fits, and ii) calculate the correction indices (see section  
166 on statistical analyses).

167 2) Multispecies mixtures, with the five selected species. We first combined them  
168 at the same proportion (20:20:20:20:20; one mixture with three replicates, N = 3) and  
169 then we created communities where the percentage of one species (either *Helianthemum*  
170 *cinereum* or *Stipa pennata*) was progressively increased and that of the other four was  
171 maintained equal (10.0:22.5:22.5:22.5:22.5, 40:15:15:15:15, 60:10:10:10:10, and  
172 80:5:5:5:5; eight types of mixtures with three replicates, N = 24). These two species  
173 were chosen because they had shown either relatively lower or higher amplification in a  
174 previous pilot study (data not shown). The multispecies mixtures were used to validate  
175 the calculated linear fit parameters and correction indices. See Table S1 for details on  
176 the composition of each type of mixture.

177

### 178 **2.3 DNA metabarcoding on root mixtures mock communities**

179 DNA was extracted from each mixture (and four isolation blanks) in the lab at  
180 Universidad Rey Juan Carlos using the DNEasy Plant Minikit (Qiagen, CA, USA) and  
181 shipped to the AllGenetics laboratories (AllGenetics & Biology SL, A Coruña, Spain).  
182 For library preparation, we amplified a fragment of the *rbcL* chloroplast gene (550 bp)  
183 using primers *rbcLa-F* (5' ATGTCACCACAAACAGAGACTAAAGC3'; Levin *et al.*  
184 2003 and *rbcLa-R* (5' GTAAAATCAAGTCCACCRCG 3'; Kress *et al.* 2009), to which  
185 the Illumina sequencing primer sequences were attached at the 5' ends. We selected the  
186 *rbcL* region because it has repeatedly been shown to be a robust barcode for plants  
187 (Hollingsworth *et al.* 2009; Kress *et al.* 2009), and because it allowed the taxonomic  
188 identification of most members of the entire study community at the species level. A  
189 series of two PCRs were carried out, the first to amplify the selected region and the



190 second to attach the index sequences required for multiplexing different libraries in the  
191 same sequencing pool. PCRs were carried out in a final volume of 25  $\mu$ l, containing 2.5  
192  $\mu$ l of template DNA, 0.5  $\mu$ M of the primers, 12.5  $\mu$ l of Supreme NZYTAq 2x Green  
193 Master Mix (NZYTech, Lisboa, Portugal), and ultrapure water up to 25  $\mu$ l. The reaction  
194 mixture consisted of an initial denaturation at 95  $^{\circ}$ C for 5 min, followed by 30 cycles of  
195 95  $^{\circ}$ C for 30 s, 52  $^{\circ}$ C for 30 s, 72  $^{\circ}$ C for 30 s, and a final extension step at 72  $^{\circ}$ C for 10  
196 minutes. The second PCR had identical conditions but only 5 cycles and 60  $^{\circ}$ C as the  
197 annealing temperature. Two negative controls with no DNA were included to check for  
198 contamination during library preparation. Portugal), The libraries were run on agarose  
199 gels stained with GreenSafe (NZYTech, Lisboa, and their size visualized under UV  
200 light. They were then purified using the Mag-Bind RXNPure Plus magnetic beads  
201 (Omega Biotek, GA, USA), pooled in equimolar amounts and sequenced in a run of the  
202 MiSeq PE300 (Illumina, CA, USA).

203         The quality of the Illumina paired-end raw data was checked using FastQC  
204 ([www.bioinformatics.babraham.ac.uk/projects/fastqc](http://www.bioinformatics.babraham.ac.uk/projects/fastqc)), and the raw reads were quality-  
205 filtered using Geneious 11.1.2 ([www.geneious.com](http://www.geneious.com)). The PCR primers were removed  
206 and a region at the 3' end of each file was trimmed according to the average Phred score  
207 (minimum Phred quality score of 20). Since the amplicons were too long to allow for  
208 the overlap of forward and reverse reads, R1 and R2 reads were concatenated using the  
209 fuse.sh script implemented in the BBmap package (Bushnell 2014). The sequences were  
210 labelled (demultiplexed) using the script multiple split libraries.py implemented in  
211 Qiime (Caporaso et al. 2010). The label is added to the headers of the FASTQ file in  
212 order to identify each sample when sequences are combined to perform downstream  
213 analysis. The resulting FASTA file was processed using the VSEARCH bioinformatic  
214 tool (Rognes et al. 2016). Sequences were dereplicated (-derep fulllength), clustered at a

215 similarity threshold of 100 % (-cluster fast,--centroids option), and sorted (-sortbysize).  
216 Artifacts (such as point mutations and chimeras) that may be generated during PCR and  
217 sequencing were filtered during the bioinformatic pipeline. *De novo* chimera detection  
218 was carried out using the UCHIME algorithm (Edgar et al. 2011) implemented in  
219 VSEARCH.

220 The taxonomic identification was performed using an in-house constructed  
221 reference database containing the representative *rbcL* sequences (553bp) of 45 species  
222 from 18 families of plants from the study community that had been collected in the  
223 same study site and individually sequenced (see sequences in XX). Since the query  
224 sequences mapped only to the 5' and 3' ends of the reference sequences, the central  
225 region of the reference sequences was previously deleted to perfectly match the query  
226 sequences, resulting in a final length of 517bp (see a similar approach in Vizcaíno *et al.*  
227 2018). The taxonomic identification was performed by querying the clustered  
228 representative sequences against our reference database using the -usearch global option  
229 of VSEARCH with a 99% similarity threshold. Finally, the script mesas-uc2clust.py  
230 was used to obtain an OTU table listing the number of sequences from each OTU found  
231 in each sample. Based on the results of this table, a quality-filtering was applied to  
232 remove the OTUs with a number of sequences lower than 0.005% of the total number of  
233 sequences (Bokulich et al. 2013; Edgar 2013).

234 In DNA metabarcoding studies, it has been observed that a low percentage of the  
235 reads of a library can be assigned to another library. This phenomenon, known as  
236 mistagging or index jumping is the result of the misassignment of the indices during  
237 library preparation, sequencing, and/or demultiplexing steps (Esling, Lejzerowicz &  
238 Pawlowski 2015; Bartram et al. 2016). To correct for this, the low abundance OTUs of  
239 each sample (0.1% threshold) were removed. Finally, only the OTUs that matched any

240 reference sequence in the database at a similarity of 99% were kept in the OTU table.  
241 The unidentified OTUs were removed from the OTU table for downstream analysis.  
242 These unidentified OTUs only accounted for an average 0.90% of the total reads before  
243 filtering.

244

## 245 **2.4 Statistical analyses**

246 All the OTUs assigned to the same species were combined before analysis. For each  
247 sample, we calculated the percentage of DNA reads assigned to each species, as the  
248 number of reads for each species divided by the total number of reads in the sample. To  
249 check whether the percentage of reads reflected the species-specific percentage of  
250 biomass in each mixture, and to improve our inference ability in the cases where there  
251 was not a reliable match between both aspects, we used two different methods: 1) linear  
252 model fits and 2) creation of correction indices (relative correction factors *sensu*  
253 Thomas et al. 2016).

254       1. *Linear model fits and match to the identity function*: for each species, we used  
255 the pairwise mixtures where it was present ( $N = 36$ ) to compute the best linear fit  
256 ( $y = ax + b$  in intercept-slope form, where  $a$  is the slope of the line and  $b$  is the  
257 intercept) between the percentage of retrieved reads ( $y$ ) and the percentage of biomass  
258 ( $x$ ; 20, 50 or 80%), using the *lm* function in R (R Core Team 2017). In order for the  
259 percentage of reads to be used directly as an estimate of the species' biomass percentage  
260 in a sample, the equation obtained for a given species would have to closely match the  
261 identity function, i.e. a linear equation where the intercept is not significantly different  
262 from zero and the slope is not significantly different from one ( $y = x$ ). If this is the  
263 case, the percentage of reads found after sequencing and filtering could be directly used  
264 to estimate the original percentage of biomass in the sample. Therefore, we used the

265 parameter estimates from the equations fitted to the data of each species to verify the  
 266 significance of the test on the intercept ( $b = 0$ ), and performed a two-tailed t-test to  
 267 compare the slope of the fit to a slope of one ( $a = 1$ )(see a similar approach on Diaz-  
 268 Real et al. 2015).

269 To assess whether the inclusion of other species affected the relationship  
 270 between the percentage of biomass of each species in the samples and the percentage of  
 271 retrieved reads, we fitted linear equations to the data including the multispecies  
 272 mixtures ( $N = 63$  for each species). Again, we tested whether the intercept was different  
 273 from zero and whether the slope of the equation was significantly different from one.  
 274 We finally compared the slope of both fits (only pairwise mixtures vs. all mixtures),  
 275 performing an analysis of covariance (ANCOVA) for each species with type of  
 276 community (pairwise vs. all) as predictor, percentage of biomass as covariate and  
 277 percentage of reads as dependent variable. A significant interaction between the  
 278 predictor and the covariate indicates that the slope fitted to the different communities is  
 279 not the same. Different intercepts and/or slopes in both equations would indicate that the  
 280 fits are affected by the species richness and composition of the samples, and therefore  
 281 that employing the percentage of sequences of DNA as estimation of biomass in  
 282 communities of different richness would be severely biased, even for the same species.

283 2. *Correction indices*: for each species, we calculated a percentage-specific  
 284 correction index using the pairwise mixtures, based on the percentage of biomass in the  
 285 sample and the percentage of reads retrieved (relative correction factors presented in  
 286 Thomas et al. 2016). Specifically, the correction index for a species  $A$  at percentage  $p$   
 287 was calculated as:

$$Correction\ index_{A_p} = \frac{Biom\ sp_{A_p}}{100 - Biom\ sp_{A_p}} * \frac{\sum_{i \neq A}^s \sum_{j=1}^r Reads\ sp_{ij}}{\sum_{i \neq A}^s \sum_{j=1}^r Reads\ sp_{Aij}}$$

288 Where  $r$  is the number of replicates for each type of community (i.e.,  $r = 3$ ),  $s$  is the  
289 number of species considered (i.e.,  $s = 5$ ),  $Biom\ sp_{A_p}$  is the percentage of biomass of  
290 the species  $A$ ,  $Reads\ sp_{A_{ij}}$  is the number of reads obtained for species  $A$  in combination  
291 with species  $i$  in each replicate  $j$ , and  $Reads\ sp_{ij}$  is the number of reads obtained for  
292 each of the other species  $i$  in replicate  $j$ . A correction index was computed for each  
293 species and percentage by averaging the number of reads in all the mixtures where it  
294 was combined with the all the other species at a specific percentage. For instance, the  
295 20% correction index for *Thymus vulgaris* was computed using all the samples where  
296 the percentage of biomass of *Thymus vulgaris* was 20% and that of the second species  
297 (either one of the other four) was 80%. When the correction index is  $\approx 1$ , the percentage  
298 of reads retrieved robustly reflects the percentage of biomass on the sample. When the  
299 correction index is well above or below 1, the percentage of reads obtained is  
300 underestimated or overestimated compared to the percentage of biomass, respectively.  
301 We computed a 20, 50, 80% correction index and an average index (average of the  
302 three) for each species.

303 The correction indices computed for each species were used to recalculate the  
304 number of reads of each species in the multispecies mixtures (by multiplying the index  
305 by the retrieved number of reads), and subsequently transform these numbers to  
306 percentage of reads. This was performed for all indices of each species (20, 50, 80% and  
307 average). Then, these corrected percentages were compared to the actual percentage of  
308 biomass of each species in the multispecies mixtures. To assess the correction ability of  
309 the computed indices, we calculated the error for each type of community (averaging  
310 the three replicates). For each species and percentage of biomass, we computed the  
311 average absolute difference between the actual percentage of biomass and the corrected  
312 percentage of reads. We also computed the average absolute differences for the original,

313 uncorrected estimates. To assess the effect of the correction and whether it was similar  
314 for all species, we used a two-way ANOVA, with correction (corrected vs. uncorrected)  
315 and species as predictors, and absolute error as the dependent variable. This was  
316 followed by species-specific one-way ANOVAs to test the effect of the correction on  
317 each species individually. Lower average error in the corrected percentages would  
318 indicate that the correction indices based on the pairwise mixtures could improve the  
319 estimates in field-collected samples of different species richness and composition.

320 An alternative method to compute correction indices is the use of a control  
321 species (see details in Thomas et al. 2016). In this approach, the correction index for  
322 each target species is computed using only the pairwise mixtures where the target and  
323 the control species are present. We computed species-specific correction indices using  
324 the 50% pairwise mixtures with three different control species: *L. suffruticosum*, *H.*  
325 *cinereum* and *S. pennata*. To assess whether the use of control-based correction indices  
326 also improved the percentages of biomass, we again used these indices to correct the  
327 percentage of reads of each species in the multispecies mixtures. Finally, for each  
328 species and multispecies mixture, we computed the average error (as defined above) of  
329 the corrections using the different control-based indices.

330

### 331 **3. RESULTS**

#### 332 **3.1 Performance of DNA metabarcoding with root mixture mock communities**

333 For most samples (> 94%), DNA metabarcoding successfully recreated the species  
334 composition of the mixtures (mock communities), i.e. all the species added to a mixture  
335 were found during sequencing. Only for a few samples where the percentage of biomass  
336 of *Stipa pennata* was low ( $\leq 20\%$ ), no sequences for this species were recovered  
337 (Supporting Information Table S1). Five species from the study plant community that  
338 were not added to the root mixtures were also detected in a few samples (N = 16), but

339 were always found at low percentages (range: 0.13-15.06%, average: 2.65%; see Table  
340 S1 for details on retrieved compositions and percentage of reads of each species and  
341 sample). Specifically, in more than 90% of the samples, more than 95% of the retrieved  
342 sequences were assigned exclusively to the species added to each mixture (Supporting  
343 Information Fig. S1).

344

### 345 **3.2 Evaluation of linear model fits and match to the identity function**

346 For the five selected species, the linear models fitted to the data were highly significant  
347 and had a positive slope, i.e. the percentage of reads of each species increased when the  
348 percentage of biomass of the species in the sample also increased (adjusted  $R^2 = 0.32$ -  
349  $0.50$ ; Fig. 2, Supporting Information Table S2). However, both the slope and the  
350 intercept of the fitted models differed among species. Specifically, the data obtained  
351 from the pairwise mixtures of *B. fruticescens*, *L. suffruticosum* and *T. vulgaris* could be  
352 fitted to a linear equation with an intercept not significantly different from zero ( $b = 0$ ,  
353 Table S2). Similarly, the t-tests showed that the slopes of the lines for these species  
354 were not significantly different from one ( $P > 0.34$  in all cases; Table S2). This  
355 indicates that the percentage of reads retrieved for these three species may be used  
356 directly to estimate the percentage of biomass on the samples. Conversely, the lines  
357 fitted to *H. cinereum* and *S. pennata* had intercepts significantly different from zero  
358 (significantly higher/lower than zero for *H. cinereum* and *S. pennata*, respectively;  
359 Table S2). Similarly, the slopes were significantly different from one for both species ( $P$   
360  $= 0.02$  and  $P < 0.001$  for *H. cinereum* and *S. pennata*, respectively), indicating that the  
361 percentage of reads was consistently higher (*H. cinereum*) or lower (*S. pennata*) than  
362 the percentage of biomass in the sample.

363           The models including all mixtures showed similar results to those fitted only  
364 with the pairwise mixtures (Table S2). Indeed, for all species except *S. pennata*, the  
365 slopes of both equations (only pairwise vs. all mixtures) were not significantly different  
366 (not significant interaction ‘percentage of biomass  $\times$  type of community’ in ANCOVA).  
367 For *B. fruticescens*, *L. suffruticosum* and *T. vulgaris*, the intercept and the slope of the  
368 lines fitted to all the mixtures were not significantly different from zero and one,  
369 respectively, matching again the identity function, and the fit improved for all species  
370 (adjusted  $R^2 = 0.56-0.66$ ; Fig. 2, Table S2). This indicates that species richness and  
371 composition did not significantly alter the fits for these three species. Conversely, when  
372 the multispecies mixtures were added to the data of the remaining species (*H. cinereum*  
373 and *S. pennata*), the equations had again intercepts significantly different from zero and  
374 slopes significantly higher/lower than one ( $P < 0.001$  for both *H. cinereum* and *S.*  
375 *pennata*; Fig. 2).

376

### 377 **3.3 Evaluation of correction indices**

378 We found a wide variation among species for the correction indices computed with all  
379 the pairwise mixtures and a control species (Fig. 3 and Supporting Information Table  
380 S3). Using all pairwise mixtures, the correction indices for *B. fruticescens*, *L.*  
381 *suffruticosum* and *T. vulgaris* were close to one, and slightly increased when the  
382 percentage of biomass in the sample increased. However, for the remaining species, the  
383 correction indices were much lower (*H. cinereum*) or much higher (*S. pennata*) than one  
384 (Fig. 3, lower panels), indicating a consistent overestimation and underestimation of the  
385 percentage of reads compared to biomass percentages.

386           When the number of reads of the multispecies mixtures were recalculated using  
387 the 50% correction indices, the recalculated percentage of reads closely matched the



388 actual biomass percentage in the multispecies samples (Fig. 4), and the average absolute  
389 error (absolute difference between the actual percentage of biomass and the percentages  
390 of reads) was significantly reduced (Fig. 5, significant differences in the average  
391 absolute error between corrected and uncorrected percentage of reads,  $P < 0.0001$ ). This  
392 error reduction was not equal among species (significant ‘species  $\times$  correction’  
393 interaction,  $P < 0.001$ ), and was especially relevant for *H. cinereum* and *S. pennata*,  
394 where their overestimation and underestimation in the uncorrected percentage of reads  
395 were significantly improved when the correction indices were applied (e.g. Fig. 4 c and  
396 d). For these species, the error between reads and biomass was significantly reduced  
397 after correction (Fig. 5). For *B. fruticescens*, *L. suffruticosum* and *T. vulgaris*, the  
398 recalculation of reads with the correction indices also improved the match between the  
399 percentage of reads and biomass (e.g. Fig. 4 g and i), although the reduction of the error  
400 was not significant for these species (Fig. 5). The use of the 50%, 80% and average  
401 correction indices rendered very similar results (from 17% error in the uncorrected  
402 samples to  $\approx 6\%$  in the corrected percentages), but the correction of the proportions was  
403 lower when the 20% correction indices were used ( $\approx 9\%$  overall error in the corrected  
404 percentages).

405 Similarly, the use of a control species to calculate the correction indices  
406 improved the estimation of percentages of biomass in most cases (Supporting  
407 Information Table S4), but the correction of the proportions varied depending on the  
408 choice of control species (Table S4), and the error reduction was on average lower than  
409 when all pairwise mixtures were used to calculate the correction indices (Supporting  
410 Information Fig. 2).

411

#### 412 **4. DISCUSSION**

413 Our study provides a straightforward and simple protocol to overcome one of the main  
414 shortcomings in DNA metabarcoding studies, the estimation of species' relative  
415 abundance based on the percentage of DNA sequences (reads) recovered. Through the  
416 use of purposefully-designed root mock communities, we test the efficacy of two  
417 complementary and easy-to-implement methods and provide robust estimates of plant  
418 biomass percentages in realistic multispecies samples. This is, to the best of our  
419 knowledge, the first study to validate a quantitative DNA metabarcoding in plant  
420 communities using root mixtures.

421 The use of metabarcoding is revolutionizing plant ecology studies, since  
422 detection of the so-called hidden diversity provides new insights to open questions in  
423 this field (see e.g. Yoccoz et al. 2012). However, the possibility of using DNA  
424 metabarcoding results to estimate species' abundances has been a subject of debate  
425 since the onset of this methodology. Due to reported inconsistencies in past attempts,  
426 recent revisions suggest that a conservative approach may be to treat metabarcoding  
427 results as presence-absence data (Deiner et al. 2017; Porter & Hajibabaei 2018).  
428 However, our study suggests that accurate quantification of species roots' biomass may  
429 be robustly done, provided that previous quantification studies using mock communities  
430 with target species are performed. Importantly, results from mock communities of root  
431 mixtures (both fitted models and corrected read percentages) can be then safely used to  
432 robustly estimate root biomass in field-collected samples, since estimated biomass  
433 percentages (using both methods) were not significantly altered by species composition,  
434 richness and species' relative abundance in the samples.

435 Our results have important implications for plant community ecology.  
436 Understanding how and to what extent stochastic and deterministic processes determine  
437 plant coexistence and community assemblages in plant communities remains an

438 unresolved question, despite the intense research effort devoted to this topic over the  
439 last decades (Götzenberger et al. 2012; Gravel et al. 2006; HilleRisLambers et al. 2012).  
440 A few authors have recognized that part of this knowledge gap could be filled if we  
441 complement our current framework, mainly based on characterization of aboveground  
442 processes, expanding our focus belowground (Bever et al. 2010; Wardle et al. 2004). In  
443 this context, DNA metabarcoding has successfully been used in a few instances to  
444 describe patterns of species richness and its distribution belowground (Hiiesalu et al.  
445 2012; Jones et al. 2011; Kesanakurti et al. 2011), but quantification attempts were  
446 lacking. The prospect of using DNA metabarcoding on root mixtures to detect not only  
447 the presence of species but also to estimate species' abundances constitutes a step  
448 further towards a deeper understanding of plant coexistence and community  
449 assemblages, especially at the fine scales where roots interact. Knowledge on the  
450 patterns of root biomass distribution will provide insights on the correspondence  
451 between above- and belowground distributions, plant-plant interactions and plant-soil  
452 feedbacks (Brandt et al. 2013; Kulmatiski et al. 2008).

453         An ideal scenario for DNA metabarcoding studies would be that the proportion  
454 of DNA sequences obtained after high-throughput sequencing closely reflected the  
455 percentage of biomass of each species in the bulk sample, irrespective of the sample  
456 composition and the relative occurrence of each species. If this was true for our plant  
457 community, the percentage of DNA sequences assigned to each species could be readily  
458 used to estimate the percentage of root biomass in field-collected samples. For the five  
459 study species, we indeed found a positive and highly significant relationship between  
460 the percentage of biomass in the pairwise mixtures and the percentage of reads  
461 recovered for each species (Fig. 2). However, the parameters of the statistical  
462 relationship (slope and intercept) widely varied among species, and for two of them (*H.*

463 *cinereum* and *S. pennata*), the best fit rendered biased estimates of biomass percentages,  
464 despite the observed positive correlation. For instance, for *H. cinereum*, the estimated  
465 percentage of biomass using the fitted equation for the pairwise communities with 20,  
466 50 or 80% biomass rendered 60, 75 and 90% biomass estimates, respectively, due to the  
467 high intercept of the fitted line (the opposite, i.e. a sharp underestimation of biomass  
468 proportions, occurred for *S. pennata*). These results point out that a significant positive  
469 relationship between percentage of biomass and percentage of reads is not sufficient to  
470 transform presence-absence data into quantitative estimates (despite its current wide  
471 use). To robustly achieve the latter, the line fitted for a given species would need to be  
472 statistically equivalent to the identity function. For three of our study species (*T.*  
473 *vulgaris*, *B. fruticescens* and *L. suffruticosum*), we found such equivalence between  
474 biomass and reads percentages. This match was not altered when the data from the  
475 pairwise mixtures was combined to the multispecies samples, suggesting that, at least  
476 for these three species, the relationship between root biomass and reads percentages is  
477 maintained regardless of the number of species (two versus five) and the species'  
478 biomass percentage (from 5 to 80%). However, even when the fit is equivalent to the  
479 identity function, the predicted abundance estimated by the linear model may be poor  
480 (e.g. fitting with a large residual error). Overall, our results call for caution on the direct  
481 use of sequence percentages to approximate relative biomass or abundance based on the  
482 existence of a positive relationship between both aspects (see e.g. Elbrecht & Leese  
483 2015; Hiiesalu et al. 2012; Pont et al. 2018) or on the mere assumption that such  
484 relationship exists (see e.g. Merges et al. 2018), and highlight the need to test the  
485 properties (statistical parameters) rather than just the existence of a significant  
486 relationship between percentage of reads and abundance for each species individually.

487 Our second approach involved the use of species-specific correction indices  
488 (based on the relative correction factors recently proposed by Thomas et al. 2016)  
489 obtained from either all the pairwise mixtures or using a control species, which were  
490 then used to correct the percentage of sequences in the multispecies samples.  
491 Importantly, these recalculated read percentages generally improved the match between  
492 the percentage of reads and the actual biomass in the multispecies mock communities  
493 (Fig. 4 and Table S4), and reduced the overall error compared to the uncorrected  
494 percentages (Figs. 5 and S2). The best results, i.e. the lowest error, was obtained when  
495 the percentage of reads were recalculated using the indices computed with all pairwise  
496 mixtures, as they closely mirrored the biomass percentages in each multispecies mock  
497 community (Fig. 4). This indicates that the use of such correction indices represents a  
498 successful way to obtain quantitative estimates in plant DNA metabarcoding studies.  
499 Several pieces of evidence support this claim. First, reliable estimates of biomass  
500 percentages were obtained for all five species after adjusting the percentages of reads,  
501 which suggests that this method can be generalized to other species in the community.  
502 This was the case even for the two species that significantly deviated from the identity  
503 function due to consistent over- and underestimation (*H. cinereum* and *S. pennata*).  
504 Indeed, the bias reduction –calculated as the difference between sequence and actual  
505 biomass percentages– between the observed and corrected percentages was higher for  
506 these two species (note that the lower bias reduction in the other three species was due  
507 to the fact that their correction indices were in all cases very close to one, i.e. no strong  
508 deviation between biomass and uncorrected reads percentages). Second, the correction  
509 index calculated for each species was computed based on pairwise mixtures of different  
510 compositions (each species combined with the other four) and then applied to  
511 multispecies samples, which highlights that these indices are robust to changes both in

512 species richness and composition. And third, the indices calculated from pairwise  
513 mixtures where the species were found at different proportions remained relatively  
514 constant (only when percentages were low, i.e. 20%, did the indices substantially  
515 differed; Fig. 3), and efficiently corrected samples where biomass percentages varied  
516 widely. This indicates that these correction indices are also relatively robust to varying  
517 species' biomass percentages. Our results concur with those by Thomas et al. (2016),  
518 the only other existing implementing correction indices, who found that control-based  
519 correction of reads proportion greatly improved relative abundances in fish mixtures.

520         This study also allowed to compare the correction ability of differently-  
521 computed indices. Although those based on all pairwise mixtures outperformed  
522 correction indices based on the use of a control species, the latter also resulted in  
523 improved biomass estimates compared to uncorrected ones. The use of a control species  
524 to compute correction indices has the advantage that the number of pairwise mixtures  
525 needed is significantly lower (e.g. in a five-species study, only four pairwise mixtures  
526 are needed if the fifth species is the control, but 10 pairwise mixtures are needed to  
527 compute indices from all pairwise mixtures), which can significantly reduce the  
528 complexity and cost of the study. However, the reduction of error widely varied  
529 depending on the choice of control species (Fig. S2), which introduces a source of  
530 uncertainty since the control species needs to be chosen a priori. In practice, the  
531 decision of how to compute correction indices will depend on a variety of factors,  
532 including the species richness of the study community, existing knowledge of the  
533 performance of species during metabarcoding, etc.

534         Our study also helped to validate the effectiveness of DNA metabarcoding using  
535 the *rbcL* region for the simultaneous identification and quantification of multiple taxa in  
536 root mixtures from Mediterranean shrublands. For most samples, the species that

537 formed each mock community were successfully recovered during sequencing. For a  
538 few samples, however, our approach recovered species –either the study species or other  
539 species from our plant community– that had not been included in those specific  
540 mixtures, although in general they accounted for a very small percentage of the DNA  
541 sequences in each sample (Fig. S1). These infrequent mismatches between the created  
542 (prepared root mixtures) and recreated (after sequencing) species composition can be  
543 due to species cross-contamination during root sampling and mock community  
544 preparation, or due to mistagging (i.e. index/tag jumping) during the DNA  
545 metabarcoding pipeline (Coissac et al. 2012; Schnell et al. 2015). Importantly, they help  
546 to identify aspects for improvement in metabarcoding studies (Deiner et al. 2017; Porter  
547 & Hajibabaei 2018). Furthermore, it is worth to note that the choice of the appropriate  
548 barcode may depend on the type of plant community and the source of DNA samples  
549 (community DNA, environmental DNA, etc.). Future studies should also incorporate  
550 several markers to determine the consistency of the correction indices across different  
551 barcodes (Hollingsworth et al. 2009).

552 In conclusion, we propose that the use of mock communities varying in species  
553 composition and biomass structure may be a useful first step for the reliable  
554 quantification of DNA metabarcoding results in other plant communities, implementing  
555 a combined approach where linear fits and correction indices are used. However, the  
556 substantial differences observed among the study species –both in the linear fits and the  
557 correction indices– indicates that quantification methods need to be applied on a  
558 species-level basis. Different sources of bias may occur during DNA extraction (e.g.  
559 differential DNA concentration per tissue biomass across species; see also Haling et al.  
560 2011) or PCR amplification (e.g. differential primer specificity; Cristescu 2014; Deiner  
561 et al. 2017; Elbrecht & Leese 2015; Pawluczyk et al. 2015; Porter & Hajibabaei 2018),

562 leading to some species being consistently under- or over-estimated during sequencing.  
563 Therefore, it is highly unlikely that biomass percentages can be estimated for all species  
564 in a community using the same linear fit or correction index. Furthermore, the ability to  
565 perform quantitative DNA metabarcoding will largely depend on the number of species  
566 in the study community, which in turn determines the amount of mock communities  
567 needed to implement corrections. In this context, prior knowledge on the species  
568 composition of the community (i.e. the existence of a robust reference library) and the  
569 selection of study species (e.g. dominant, keystone species) are critical for the  
570 successful implementation of reliable quantification methods. Finally, our results also  
571 suggest that the indiscriminate use of uncorrected percentages of sequences as a proxy  
572 for species' biomass without previous quantification tests such as the one presented here  
573 may render strongly biased results for many species.

574

#### 575 **ACKNOWLEDGEMENTS**

576 We are in debt to Joaquín Vierna, Neus Marí-Mena and Antón Vizcaino (AllGenetics)  
577 for their dedication and support during DNA metabarcoding and bioinformatic analysis.  
578 We are also grateful to Carlos Díaz (URJC) for his help in the collection of plant  
579 material. We are also grateful to three anonymous reviewers and editor for their  
580 thorough revision of our manuscript. Funding was provided by MINECO grant ROOTS  
581 (CGL2015-66809-P).

582

#### 583 **AUTHOR CONTRIBUTION**

584 AE, DSP and SM conceived and designed the study. DSP, BP and SM prepared the  
585 mock communities. JC-B prepared the reference database. AI, MC, DSP and SM



586 analyzed the data. All authors contributed to the discussion and interpretation of the  
587 results. SM wrote the manuscript, with input from all other authors.

588

## 589 DATA ACCESSIBILITY

590 Data has been deposited in Dryad (doi:10.5061/dryad.dm4t39t).

591

592

## 593 REFERENCES

- 594 Amend AS, Seifert KA, Bruns TD (2010) Quantifying microbial communities with 454  
595 pyrosequencing: does read abundance count? *Molecular Ecology* **19**, 5555-5565.
- 596 Armas C, Pugnaire FI (2011) Belowground zone of influence in a tussock grass species.  
597 *Acta Oecologica-International Journal of Ecology* **37**, 284-289.
- 598 Baird, D. J., & Hajibabaei, M. (2012). Biomonitoring 2.0: a new paradigm in ecosystem  
599 assessment made possible by next generation DNA sequencing. *Molecular*  
600 *Ecology*, 21(8), 2039-2044.
- 601 Bardgett RD, Mommer L, De Vries FT (2014) Going underground: root traits as drivers  
602 of ecosystem processes. *Trends in Ecology & Evolution* **29**, 692-699.
- 603 Bell KL, Burgess KS, Botsch JC, et al. (2019) Quantitative and qualitative assessment  
604 of pollen DNA metabarcoding using constructed species mixtures. *Molecular*  
605 *Ecology* **28**, 431-455.
- 606 Bever JD, Dickie IA, Facelli E, et al. (2010) Rooting theories of plant community  
607 ecology in microbial interactions. *Trends in Ecology & Evolution* **25**, 468-478.
- 608 Bokulich NA, Subramanian S, Faith JJ, et al. (2013) Quality-filtering vastly improves  
609 diversity estimates from Illumina amplicon sequencing. *Nature Methods* **10**, 57.
- 610 Brandt AJ, Kroon H, Reynolds HL, Burns JH (2013) Soil heterogeneity generated by  
611 plant-soil feedbacks has implications for species recruitment and coexistence.  
612 *Journal of Ecology* **101**, 277-286.
- 613 Bushnell, B. BBMap: a fast, accurate, splice-aware aligner. Lawrence Berkeley  
614 National Lab.(LBNL), Berkeley, CA (United States), 2014.
- 615 Caporaso, J. G., Kuczynski, J., Stombaugh, J., et al. (2010). QIIME allows analysis of  
616 high-throughput community sequencing data. *Nature Methods* **7**, 335.
- 617 Casper BB, Jackson RB (1997) Plant competition underground. *Annual Review of*  
618 *Ecology and Systematics* **28**, 545-570.
- 619 Coissac E, Riaz T, Puillandre N (2012) Bioinformatic challenges for DNA  
620 metabarcoding of plants and animals. *Molecular Ecology* **21**, 1834-1847.
- 621 Cristescu ME (2014) From barcoding single individuals to metabarcoding biological  
622 communities: towards an integrative approach to the study of global  
623 biodiversity. *Trends in Ecology & Evolution* **29**, 566-571.
- 624 Chacón-Labela J, Cruz M, Escudero A (2017) Evidence for a stochastic geometry of  
625 biodiversity: the effects of species abundance, richness and intraspecific  
626 clustering. *Journal of Ecology* **105**, 382-390.
- 627 Chacón-Labela J, de la Cruz M, Escudero A (2016) Beyond the classical nurse species  
628 effect: diversity assembly in a Mediterranean ~~semi~~ dwarf shrubland.  
629 *Journal of Vegetation Science* **27**, 80-88.

- 630 Chase JM, Leibold MA (2003) *Ecological niches: linking classical and contemporary*  
631 *approaches* University of Chicago Press.
- 632 Deagle BE, Thomas AC, Shaffer AK, Trites AW, Jarman SN (2013) Quantifying  
633 sequence proportions in a DNA-based diet study using Ion Torrent amplicon  
634 sequencing: which counts count? *Molecular Ecology Resources* **13**, 620-633.
- 635 Deiner K, Bik HM, Mächler E, *et al.* (2017) Environmental DNA metabarcoding:  
636 transforming how we survey animal and plant communities. *Molecular Ecology*  
637 **26**, 5872-5895.
- 638 Diaz-Real J, Serrano D, Piriz A, Jovani R (2015) NGS metabarcoding proves successful  
639 for quantitative assessment of symbiont abundance: the case of feather mites on  
640 birds. *Experimental and Applied Acarology* **67**, 209-218.
- 641 Edgar RC (2013) UPARSE: highly accurate OTU sequences from microbial amplicon  
642 reads. *Nature Methods* **10**, 996.
- 643 Edgar, R. C., Haas, B. J., Clemente, J. C., Quince, C., & Knight, R. (2011). UCHIME  
644 improves sensitivity and speed of chimera detection. *Bioinformatics*, 27(16),  
645 2194-2200.
- 646 Elbrecht V, Leese F (2015) Can DNA-based ecosystem assessments quantify species  
647 abundance? Testing primer bias and biomass—sequence relationships with an  
648 innovative metabarcoding protocol. *PLoS ONE* **10**, e0130324.
- 649 Evans NT, Olds BP, Renshaw MA, *et al.* (2016) Quantification of mesocosm fish and  
650 amphibian species diversity via environmental DNA metabarcoding. *Molecular*  
651 *Ecology Resources* **16**, 29-41.
- 652 Fahner NA, Shokralla S, Baird DJ, Hajibabaei M (2016) Large-Scale Monitoring of  
653 Plants through Environmental DNA Metabarcoding of Soil: Recovery,  
654 Resolution, and Annotation of Four DNA Markers. *PLoS ONE* 11, e0157505.
- 655 Fonseca VG (2018) “Pitfalls in relative abundance estimation using eDNA  
656 metabarcoding”. *Molecular Ecology Resources* **18**, 923-926.
- 657 Götzenberger L, de Bello F, Bråthen KA, *et al.* (2012) Ecological assembly rules in  
658 plant communities—approaches, patterns and prospects. *Biological Reviews* **87**,  
659 111-127.
- 660 Gravel D, Canham CD, Beaudet M, Messier C (2006) Reconciling niche and neutrality:  
661 the continuum hypothesis. *Ecology Letters* **9**, 399-409.
- 662 Haling RE, Simpson RJ, McKay AC, *et al.* (2011) Direct measurement of roots in soil  
663 for single and mixed species using a quantitative DNA-based method. *Plant and*  
664 *Soil* **348**, 123-137.
- 665 Harvey JB, Johnson SB, Fisher JL, Peterson WT, Vrijenhoek RC (2017) Comparison of  
666 morphological and next generation DNA sequencing methods for assessing  
667 zooplankton assemblages. *Journal of Experimental Marine Biology and Ecology*  
668 **487**, 113-126.
- 669 Hawkins J, de Vere N, Griffith A, *et al.* (2015) Using DNA metabarcoding to identify  
670 the floral composition of honey: A new tool for investigating honey bee foraging  
671 preferences. *PLoS ONE* **10**, e0134735.
- 672 Hiiesalu I, OePIK M, Metsis M, *et al.* (2012) Plant species richness belowground:  
673 higher richness and new patterns revealed by-gene-sequencing.  
674 *Molecular Ecology* **21**, 2004-2016.
- 675 Hilbert DW, Canadell J (1995) Biomass partitioning and resource allocation of plants  
676 from Mediterranean-type ecosystems: possible responses to elevated  
677 atmospheric CO<sub>2</sub>. In: *Global Change and Mediterranean-Type Ecosystems*, pp.  
678 76-101. Springer.

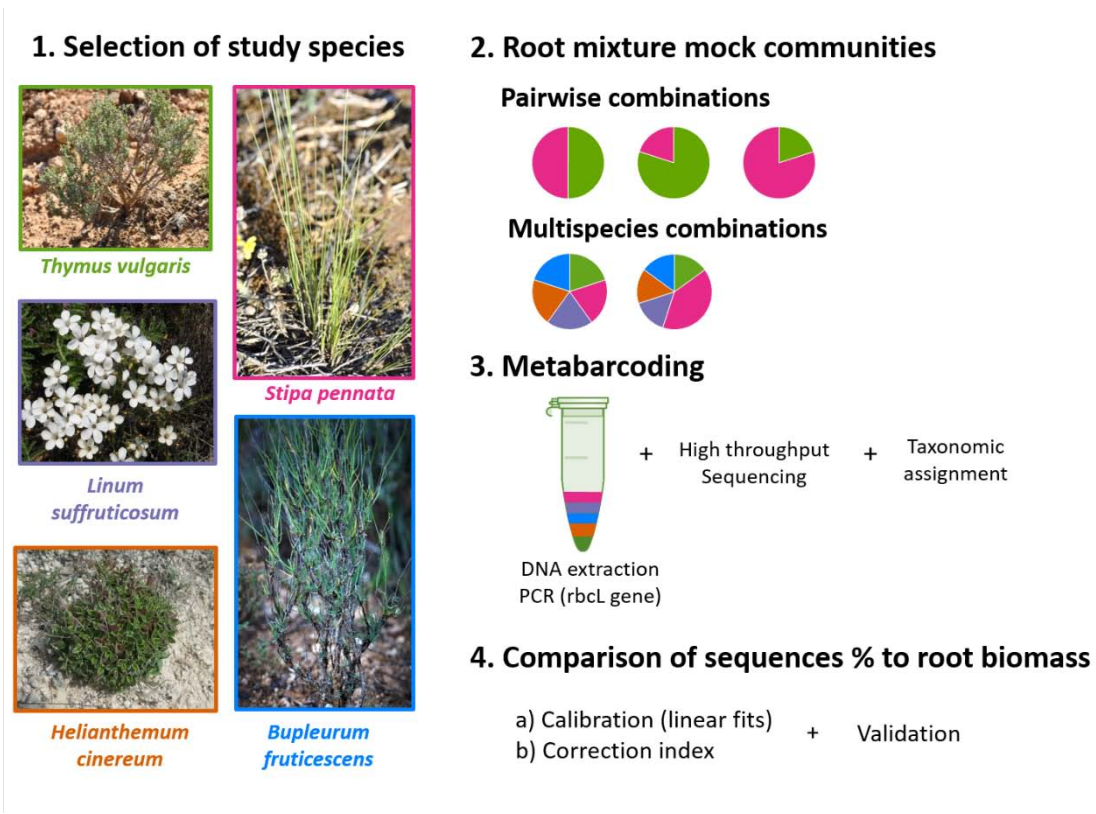
- 679 HilleRisLambers J, Adler P, Harpole W, Levine J, Mayfield M (2012) Rethinking  
680 community assembly through the lens of coexistence theory. *Annual Review of*  
681 *Ecology, Evolution, and Systematics* **43**.
- 682 Hollingsworth PM, Forrest LL, Spouge JL, *et al.* (2009) A DNA barcode for land  
683 plants. *Proceedings of the National Academy of Sciences* **106**, 12794-12797.
- 684 Jones FA, Erickson DL, Bernal MA, *et al.* (2011) The Roots of Diversity: Below  
685 Ground Species Richness and Rooting Distributions in a Tropical Forest  
686 Revealed by DNA Barcodes and Inverse Modeling. *PLoS ONE* **6**.
- 687 Jørgensen T, Kjaer KH, Haile J, *et al.* (2012) Islands in the ice: detecting past  
688 vegetation on Greenlandic nunataks using historical records and sedimentary  
689 ancient DNA Meta-barcoding. *Molecular Ecology* **21**, 1980-1988.
- 690 Kesanakurti PR, Fazekas AJ, Burgess KS, *et al.* (2011) Spatial patterns of plant  
691 diversity below-ground as revealed by DNA barcoding. *Molecular Ecology* **20**,  
692 1289-1302.
- 693 Kress WJ, Erickson DL, Jones FA, *et al.* (2009) Plant DNA barcodes and a community  
694 phylogeny of a tropical forest dynamics plot in Panama. *Proceedings of the*  
695 *National Academy of Sciences*, pnas. 0909820106.
- 696 Kress WJ, Wurdack KJ, Zimmer EA, Weigt LA, Janzen DH (2005) Use of DNA  
697 barcodes to identify flowering plants. *Proceedings of the National Academy of*  
698 *Sciences* **102**, 8369-8374.
- 699 Kulmatiski A, Beard KH, Stevens JR, Cobbold SM (2008) Plant–soil feedbacks: a  
700 meta-analytical review. *Ecology Letters* **11**, 980-992.
- 701 Levin RA, Wagner WL, Hoch PC, *et al.* (2003) Family relationships of  
702 Onagraceae based on chloroplast rbcL and ndhF data. *American Journal of*  
703 *Botany* **90**, 107-115.
- 704 Lim NKM, Tay YC, Srivathsan A, *et al.* (2016) Next-generation freshwater  
705 bioassessment: eDNA metabarcoding with a conserved metazoan primer reveals  
706 species-rich and reservoir-specific communities. *Royal Society Open Science* **3**.
- 707 McGill BJ, Etienne RS, Gray JS, *et al.* (2007) Species abundance distributions: moving  
708 beyond single prediction theories to integration within an ecological framework.  
709 *Ecology Letters* **10**, 995-1015.
- 710 Merges D, Bálint M, Schmitt I, Böhringer K, Neuschulz EL (2018) Spatial  
711 patterns of pathogenic and mutualistic fungi across the elevational range of a  
712 host plant. *Journal of Ecology* **106**, 1545-1557.
- 713 Nichols RV, Vollmers C, Newsom LA, *et al.* (2018) Minimizing polymerase biases in  
714 metabarcoding. *Molecular Ecology Resources* **18**, 927-939.
- 715 Pawlak AR, Mack RN, Busch JW, Novak SJ (2015) Invasion of *Bromus tectorum* (L.)  
716 into California and the American Southwest: rapid, multi-directional and  
717 genetically diverse. *Biological Invasions* **17**, 287-306.
- 718 Pawluczyk M, Weiss J, Links MG, *et al.* (2015) Quantitative evaluation of bias in PCR  
719 amplification and next-generation sequencing derived from metabarcoding  
720 samples. *Analytical and Bioanalytical Chemistry* **407**, 1841-1848.
- 721 Philippot L, Raaijmakers JM, Lemanceau P, van der Putten WH (2013) Going back to  
722 the roots: the microbial ecology of the rhizosphere. *Nature Reviews*  
723 *Microbiology* **11**, 789-799.
- 724 Piñol J, Mir G, Gomez-Pompa P, Agustí N (2015) Universal and blocking primer  
725 mismatches limit the use of high-throughput DNA sequencing for the  
726 quantitative metabarcoding of arthropods. *Molecular Ecology Resources* **15**,  
727 819-830.

- 728 Polz MF, Cavanaugh CM (1998) Bias in template-to-product ratios in multitemplate  
729 PCR. *Applied and Environmental Microbiology* 64, 3724-3730.
- 730 Pont D, Rocle M, Valentini A, *et al.* (2018) Environmental DNA reveals quantitative  
731 patterns of fish biodiversity in large rivers despite its downstream transportation.  
732 *Scientific reports* 8, 10361.
- 733 Poorter H, Niklas KJ, Reich PB, *et al.* (2012) Biomass allocation to leaves, stems and  
734 roots: meta-analyses of interspecific variation and environmental control. *New*  
735 *Phytologist* 193, 30-50.
- 736 Porter TM, Hajibabaei M (2018) Scaling up: A guide to **high**throughput genomic  
737 approaches for biodiversity analysis. *Molecular Ecology* 27, 313-338.
- 738 Porter TM, Shokralla S, Baird D, Golding GB, Hajibabaei M (2016) Ribosomal DNA  
739 and Plastid Markers Used to Sample Fungal and Plant Communities from  
740 Wetland Soils Reveals Complementary Biotas. *PLoS ONE* 11, e0142759.
- 741 Rognes, T., Flouri, T., Nichols, B., Quince, C., & Mahé, F. (2016). VSEARCH: a  
742 versatile open source tool for metagenomics. *PeerJ*, 4, e2584.
- 743 Schenk HJ, Jackson RB (2002) Rooting depths, lateral root spreads and below-  
744 ground/above-ground allometries of plants in water-limited ecosystems. *Journal*  
745 *of Ecology* 90, 480-494.
- 746 Schnell IB, Bohmann K, Gilbert MTP (2015) Tag jumps illuminated—reducing  
747 sequence-to-sample misidentifications in metabarcoding studies. *Molecular*  
748 *Ecology Resources* 15, 1289-1303.
- 749 Schuster SC (2007) Next-generation sequencing transforms today's biology. *Nature*  
750 *Methods* 5, 16.
- 751 Silva JS, Rego FC (2003) Root distribution of a Mediterranean shrubland in Portugal.  
752 *Plant and Soil* 255, 529-540.
- 753 Taberlet P, Coissac E, Pompanon F, Brochmann C, Willerslev E (2012) Towards  
754 next-generation biodiversity assessment using DNA metabarcoding. *Molecular*  
755 *Ecology* 21, 2045-2050.
- 756 Thomas AC, Deagle BE, Eveson JP, Harsch CH, Trites AW (2016) Quantitative DNA  
757 metabarcoding: improved estimates of species proportional biomass using  
758 correction factors derived from control material. *Molecular Ecology Resources*  
759 16, 714-726.
- 760 Vellend M (2010) Conceptual synthesis in community ecology. *The Quarterly Review*  
761 *of Biology* 85, 183-206.
- 762 Wardle DA, Bardgett RD, Klironomos JN, *et al.* (2004) Ecological linkages between  
763 aboveground and belowground biota. *Science* 304, 1629-1633.
- 764 Yoccoz N, Bråthen K, Gielly L, *et al.* (2012) DNA from soil mirrors plant taxonomic  
765 and growth form diversity. *Molecular Ecology* 21, 3647-3655.

766

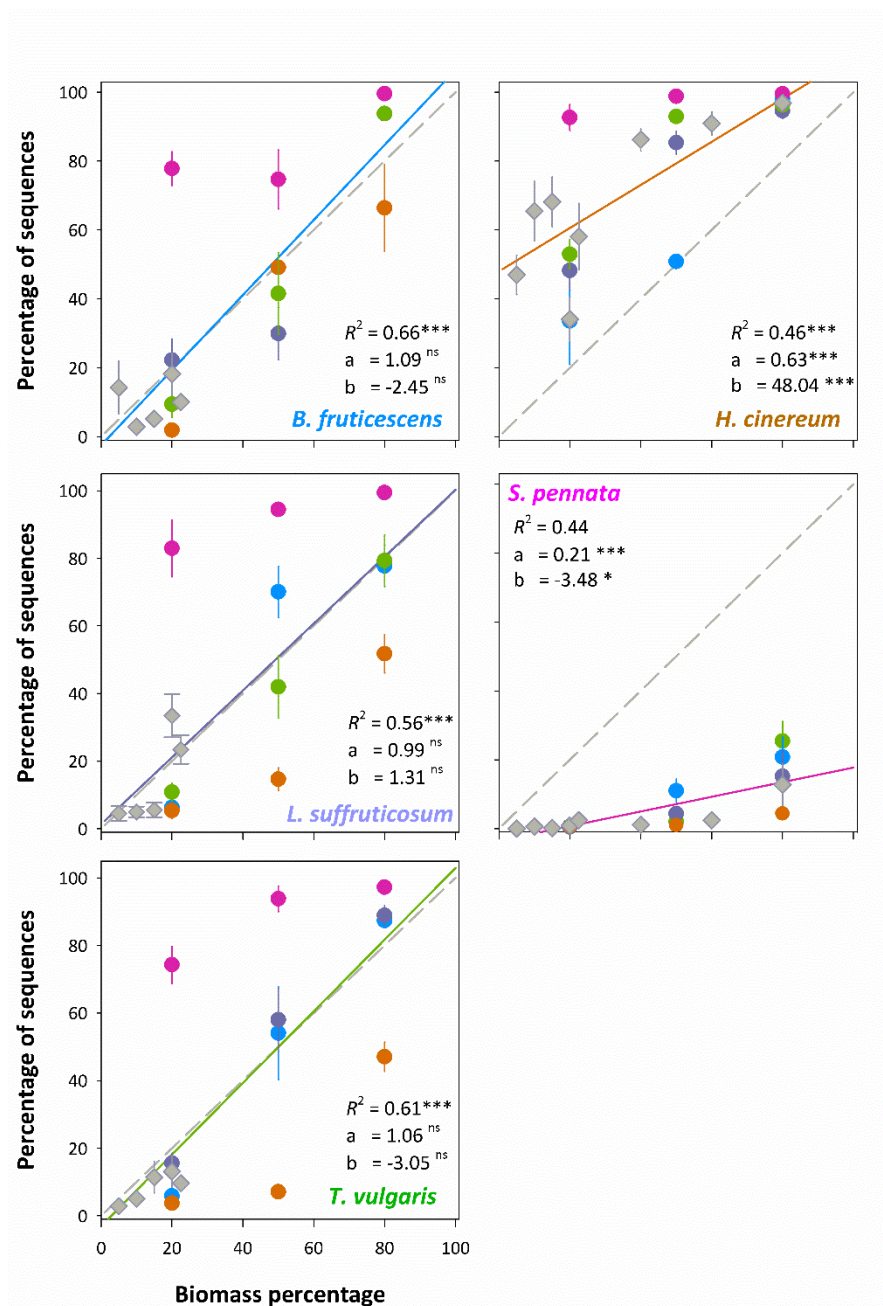
767

768 **Figure 1.** Workflow implemented to validate DNA metabarcoding quantification  
 769 methods. 1) Selection of the most dominant species in the study plant community  
 770 (based on the number of individuals). 2) Creation of the pairwise and multispecies  
 771 mixtures, i.e. root mock communities of known composition and varying percentage of  
 772 biomass proportion (see Table S1 for detailed information on composition of all mock  
 773 communities). 3) DNA metabarcoding and bioinformatics pipeline: DNA extraction,  
 774 PCR (*rbcL* gene), next-generation-sequencing and taxonomic assignment. 4)  
 775 Quantification methods and validation: calculation of the percentage of reads (DNA  
 776 sequences) assigned to each species in each mixture and comparison to the actual  
 777 percentage of biomass in the sample via linear fits and correction indices. Validation of  
 778 both methods with multispecies mixtures.  
 779  
 780



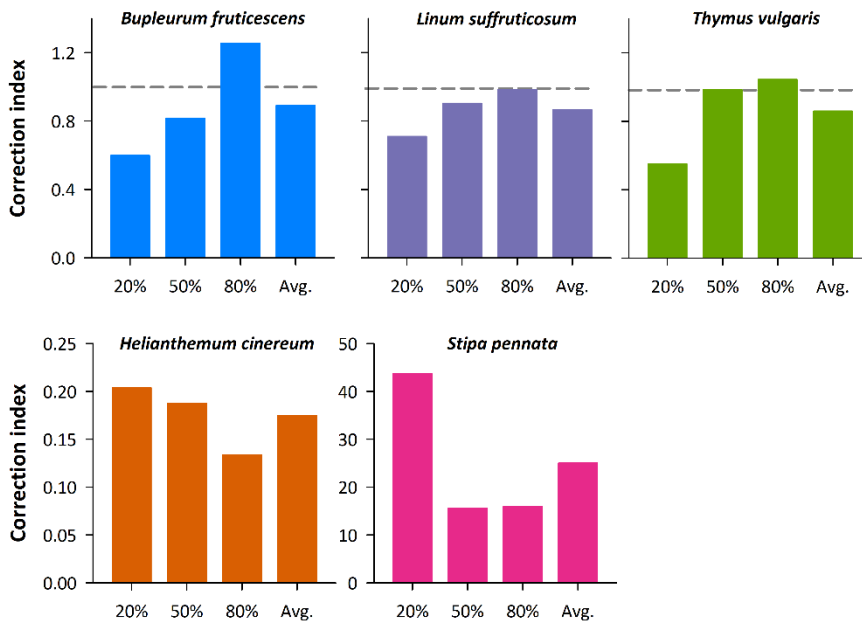
781

782 **Figure 2.** Percentage of DNA reads recovered as a function of the percentage of  
 783 biomass in each mixture (mock community). Each panel presents all the mixtures where  
 784 each species was present. Mean  $\pm$  standard error of three replicates of each mixture are  
 785 shown. Different colors of the symbols reflect the accompanying species in the pairwise  
 786 mixtures (i.e. they correspond to colors of scientific names in each panel; e.g. pink  
 787 circles in the top left panel represent pairwise mixtures of *B. fruticescens* with *S.*  
 788 *pennata*). Grey diamond shapes represent the proportion of DNA reads obtained for  
 789 each species in the multispecies mixtures. The best linear fit (colored line), adjusted  $R^2$ ,  
 790 intercept ( $b$ ) and slope ( $a$ ) estimates are also shown for each species (pairwise and  
 791 multispecies samples combined). A significant intercept and/or slope indicate  
 792 significant differences from zero and one, respectively. The grey dashed line represents  
 793 the intercept = 0 and slope = 1 fit (identity function). \*\*\*  $P < 0.001$ ; \*\*  $P < 0.01$ ; \*  $P <$   
 794  $0.05$ ; ns, not significant.  
 795



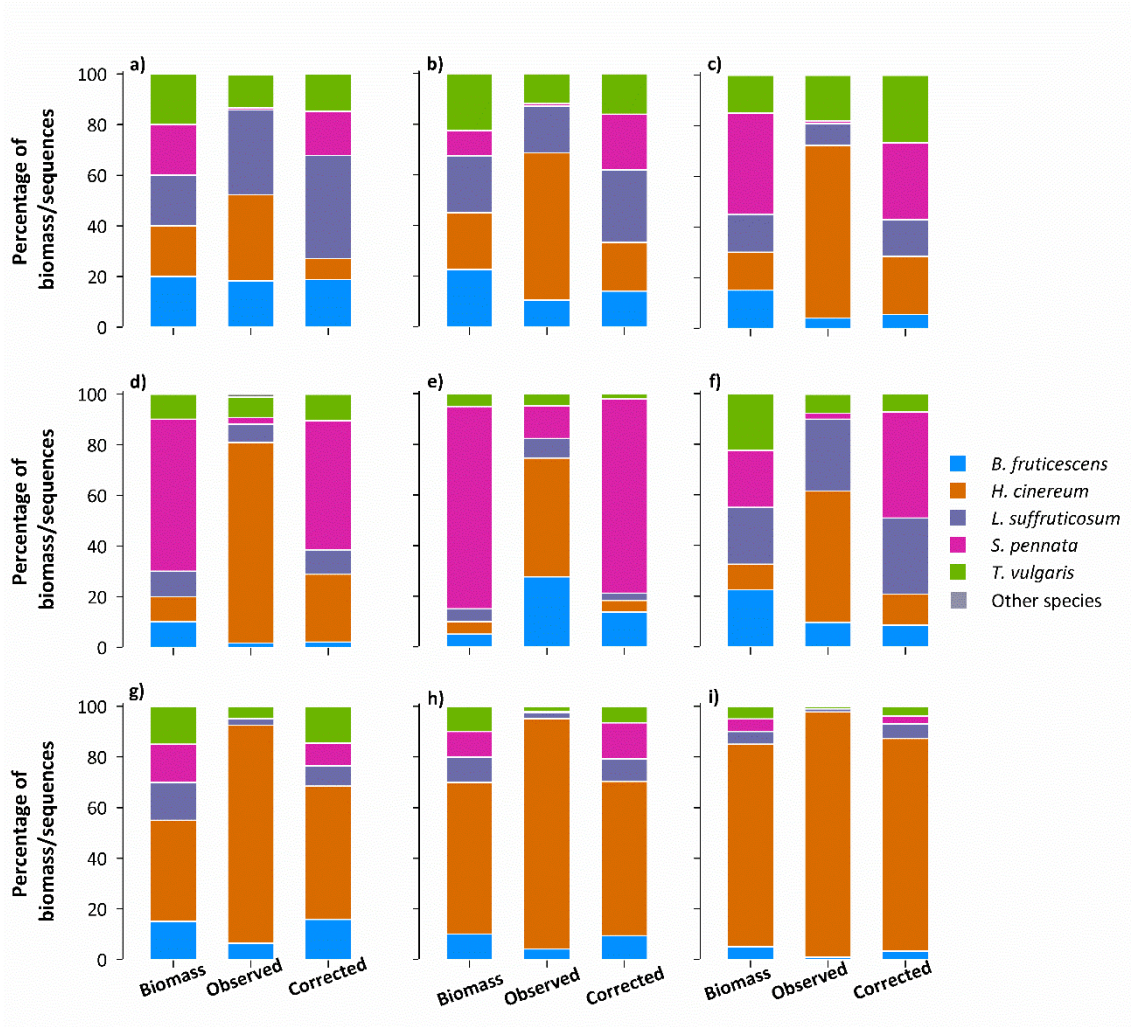
796

797 **Figure 3.** Correction indices for each species calculated from the pairwise mixtures,  
798 based on the number of reads retrieved after sequencing and the percentage of biomass  
799 proportion in each mixture. For each species, a correction index was calculated using  
800 the mixtures where the species was at 20, 50 and 80% of biomass. Avg. refers to the  
801 average correction index. The dashed grey line in the top panels represents a correction  
802 index = 1. Note that the Y axes for the species in the top panels is the same.  
803  
804



805  
806

807 **Figure 4.** Comparison of the percentage of biomass of each species in the multispecies  
 808 mixtures (left bars) to the observed (uncorrected; central bars) and the corrected (after  
 809 recalculation of the number of reads with the 50% correction indices; right bars)  
 810 percentage of reads. Panels a) to i) show a specific type of multispecies mixture  
 811 (defined by the left column), and each color represents the percentage of biomass/reads  
 812 proportion of each species, averaged for the three replicate samples of each mixture.  
 813 Gray stacks in the uncorrected bars represent the proportion of sequences retrieved from  
 814 species other than those included in the mixtures.  
 815

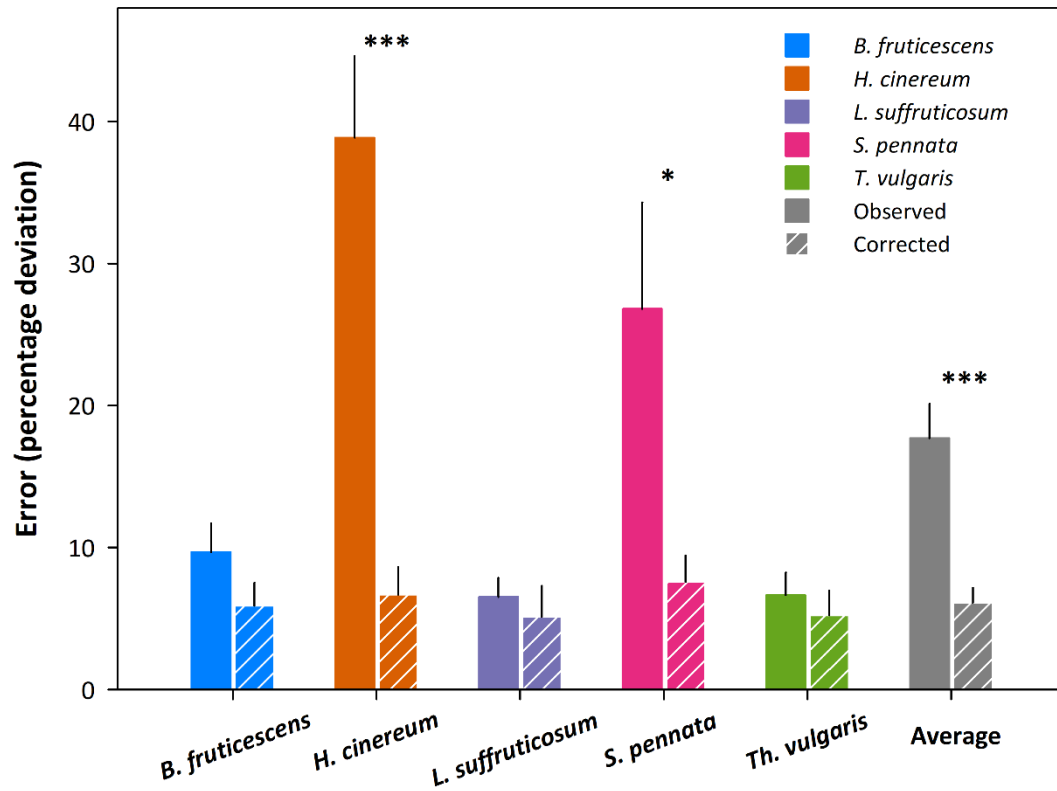


816

817



818 **Figure 5.** Average error (average absolute difference between the real percentage of  
 819 biomass and the corrected/uncorrected percentages of reads) for each species. Filled  
 820 bars: error for uncorrected reads. Striped bars: error for corrected reads. Lines represent  
 821 1 s.e. Significant differences between corrected and uncorrected deviations (one-way  
 822 ANOVA) are indicated by asterisks (\*\* $P < 0.01$ ; \* $P < 0.05$ ).



823

824

825

# MOLECULAR ECOLOGY RESOURCES

Supplemental Information for:

## Estimating belowground plant abundance with DNA metabarcoding

Silvia Matesanz, David S. Pescador, Beatriz Pías, Ana M. Sánchez, Julia Chacón-Labela, Angela Illuminati, Marcelino de la Cruz, Jesús López-Angulo, Adrián Escudero

### Table of Contents:

<b>Table S1</b>	Page 2
<b>Figure S1</b>	Page 6
<b>Table S2</b>	Page 7
<b>Table S3</b>	Page 8
<b>Table S4</b>	Page 9
<b>Figure S2</b>	Page 12

\*\*

### Instructions for Authors:

1. use of this branded Supplemental Information template is recommended, but not mandatory
2. consolidate your Supplemental files into as few documents as possible
3. if your file of Supplemental Information is very large, create a Table of Contents
4. Your Supplemental Information will not be copy-edited. Do not leave in track-changes and other editing marks. The document will be posted "as is."
5. Save as PDF if possible

# MOLECULAR ECOLOGY RESOURCES

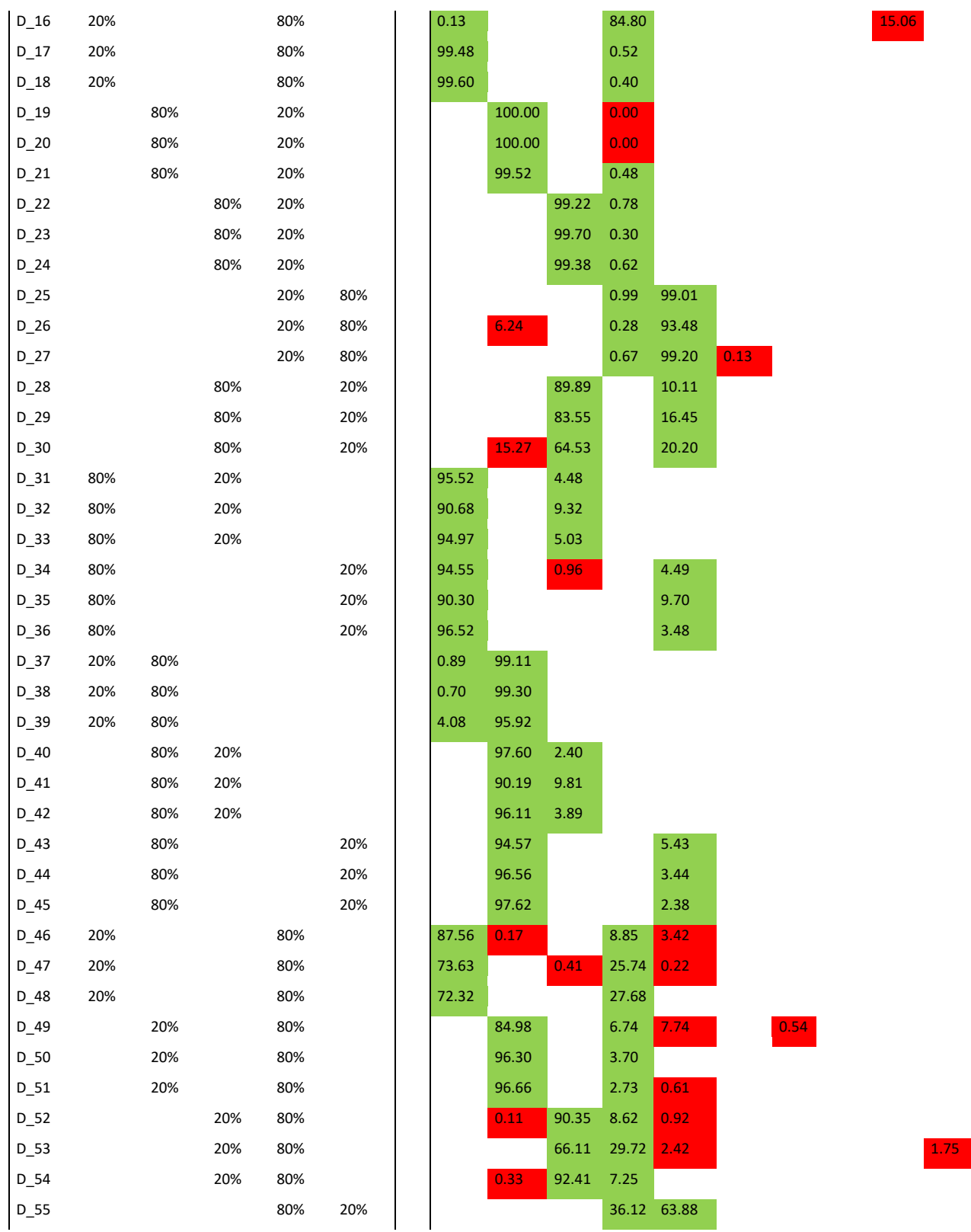
**Table S1.** Species composition and biomass proportion of each mock community (left). The right section of the table shows the positive identification of the species after sequencing (in green), and the detection of species during sequencing that were not added in the community (in red). Comm.: Community; B\_s: *Bupleurum frutescens*; H\_c: *Helianthemum cinereum*; L\_s: *Linum suffruticosum*; S\_p: *Stipa pennata*; T\_v: *Thymus vulgaris*; S\_l: *Salvia lavandulifolia*; C\_m: *Coris monspeliensis*; Q\_sp: *Quercus sp*; T\_d: *Thesium divaricatum*; L\_c: *Leuzea conifera*.

Sample	B_s	H_c	L_s	S_p	T_v	B_s	H_c	L_s	S_p	T_v	S_l	C_m	Q_sp	T_d	L_c
A_01	50%		50%			38.50		61.50							
A_02	50%		50%			14.80		85.20							
A_03	50%		50%			36.50		63.50							
A_04	50%				50%	17.70		1.00		81.30					
A_05	50%				50%	53.22	1.81	1.31		37.61				6.05	
A_06	50%				50%	53.80		3.00		43.30					
A_07	50%	50%				45.40	54.60								
A_08	50%	50%				49.80	50.20								
A_09	50%	50%				52.20	47.80								
A_10		50%	50%				89.01	10.99							
A_11		50%	50%				78.63	21.37							
A_12		50%	50%				88.38	11.62							
A_13		50%			50%		94.54			5.46					
A_14		50%			50%		92.57			7.43					
A_15		50%			50%		91.65			8.35					
A_16	50%			50%		83.41	0.23	10.34	5.41	0.28					0.33
A_17	50%			50%		83.05		0.45	10.44				3.58	2.49	
A_18	50%			50%		57.55		19.48	17.12	0.22					5.62
A_19		50%		50%			99.00		0.57	0.43					
A_20		50%		50%			98.55		1.45						
A_21		50%		50%			100.00		0.00						
A_22			50%	50%					95.54	3.99	0.47				
A_23			50%	50%					94.56	3.50	1.94				
A_24			50%	50%					93.25	5.62	1.12				
A_25				50%	50%		0.86		3.37	95.77					
A_26				50%	50%				0.87	99.13					
A_27				50%	50%		11.73		1.77	86.49					
A_28			50%		50%			25.05		74.95					
A_29			50%		50%			56.79		43.21					
A_30			50%		50%			44.12		55.88					
B_01	20%	20%	20%	20%	20%	24.47	26.03	37.22	1.09	10.90	0.29				
B_02	20%	20%	20%	20%	20%	6.08	29.47	42.08	0.51	21.85					

# MOLECULAR ECOLOGY RESOURCES

B_03	20%	20%	20%	20%	20%	24.06	46.81	20.95	1.07	6.53	0.57
C_01	22.50%	22.50%	22.50%	10%	22.50%	8.50	71.34	13.49	0.37	6.31	
C_02	22.50%	22.50%	22.50%	10%	22.50%	12.01	39.26	27.45	1.21	20.07	
C_03	22.50%	22.50%	22.50%	10%	22.50%	11.25	63.55	14.63	1.22	9.34	
C_04	15%	15%	15%	40%	15%	5.36	67.53	15.80	1.30	10.01	
C_05	15%	15%	15%	40%	15%	2.06	80.96	5.48	0.94	10.56	
C_06	15%	15%	15%	40%	15%	4.59	56.00	4.69	1.16	33.56	
C_07	10%	10%	10%	60%	10%	1.11	80.74	10.42	3.45	4.27	
C_08	10%	10%	10%	60%	10%	2.96	64.82	9.47	3.24	16.01	3.50
C_09	10%	10%	10%	60%	10%	0.94	92.01	2.18	0.74	4.13	
C_10	5%	5%	5%	80%	5%	8.36	45.86	15.03	24.78	5.98	
C_11	5%	5%	5%	80%	5%	45.49	37.77	4.31	6.64	5.79	
C_12	5%	5%	5%	80%	5%	29.00	57.31	3.93	7.03	2.74	
C_13	22.50%	10%	22.50%	22.50%	22.50%	5.80	73.96	15.00	1.70	3.54	
C_14	22.50%	10%	22.50%	22.50%	22.50%	11.34	42.99	35.24	3.34	6.93	0.16
C_15	22.50%	10%	22.50%	22.50%	22.50%	11.63	38.64	34.84	2.33	11.87	0.70
C_16	40%	15%	15%	15%	15%	16.22	80.97	1.19	0.17	1.45	
C_17	40%	15%	15%	15%	15%	1.92	85.79	4.64	0.00	7.65	
C_18	40%	15%	15%	15%	15%	1.01	91.79	1.50	0.37	5.33	
C_19	60%	10%	10%	10%	10%	0.66	95.64	1.97	0.00	1.73	
C_20	60%	10%	10%	10%	10%	11.02	84.38	1.79	0.63	2.18	
C_21	60%	10%	10%	10%	10%	0.78	92.76	3.84	0.33	2.30	
C_22	80%	5%	5%	5%	5%	0.95	95.81	1.58	0.00	1.66	
C_23	80%	5%	5%	5%	5%	1.58	96.05	1.45	0.15	0.77	
C_24	80%	5%	5%	5%	5%	0.19	98.69	0.73	0.00	0.39	
D_01	20%		80%			31.85		68.15			
D_02	20%		80%			24.08		75.92			
D_03	20%		80%			10.81		89.19			
D_04	20%			80%		15.67				84.33	
D_05	20%			80%		2.52		7.76		89.72	
D_06	20%			80%		10.09	0.25	1.65		88.00	
D_07	20%	80%				43.58	56.42				
D_08	20%	80%				87.10	12.90				
D_09	20%	80%				68.52	31.48				
D_10		20%	80%				36.90	63.10			
D_11		20%	80%				54.20	45.80			
D_12		20%	80%				53.61	46.39			
D_13		20%		80%			61.34			38.66	
D_14		20%		80%			46.93			53.07	
D_15		20%		80%			50.61			49.39	

# MOLECULAR ECOLOGY RESOURCES

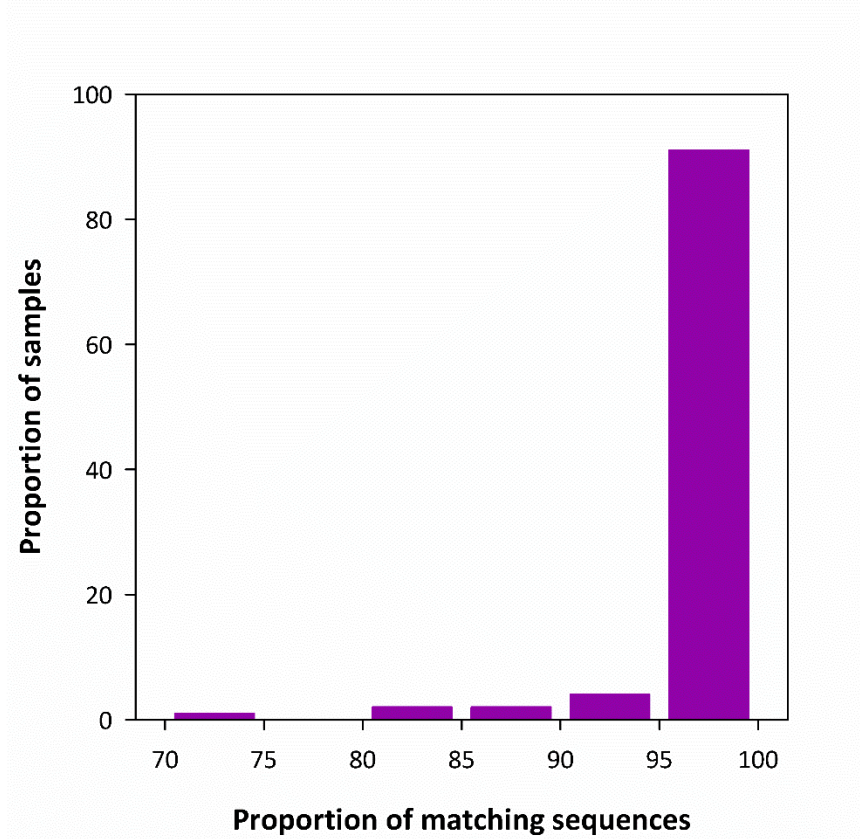


# MOLECULAR ECOLOGY RESOURCES

D_56	80%	20%	0.94	16.62	82.44
D_57	80%	20%		23.54	76.46
D_58	20%	80%	0.94	15.81	83.26
D_59	20%	80%		9.79	90.21
D_60	20%	80%		6.85	93.15

# MOLECULAR ECOLOGY RESOURCES

**Figure S1.** Proportion of samples according to the proportion of sequences assigned to the species added in each mock community.



# MOLECULAR ECOLOGY RESOURCES

**Table S2.** Intercept, slope, p-values for t-tests (hypotheses: intercept = 0 and slope = 1) and adjusted  $R^2$  of the fit between the percentage of biomass in the samples and the percentage of reads (DNA sequences).

Pairwise samples	Intercept	P-value		Adj. $R^2$	
		Intercept test	Slope		
<i>Bupleurum fruticosens</i>	5.224	0.576	0.9873	0.501	
<i>Helianthemum cinereum</i>	44.6472	<b>&lt;0.00001</b>	0.666924	<b>0.020</b>	
<i>Linum suffruticosum</i>	10.6417	0.346	0.8453	0.325	
<i>Stipa pennata</i>	-6.32159	<b>0.039</b>	0.26837	<b>&lt;0.00001</b>	
<i>Thymus vulgaris</i>	6.6994	0.543	0.9209	0.377	
<hr/>					
All samples					
<i>Bupleurum fruticosens</i>	-2.44947	0.566	1.08932	0.941	0.660
<i>Helianthemum cinereum</i>	48.1151	<b>&lt;0.00001</b>	0.6336	<b>0.012</b>	0.470
<i>Linum suffruticosum</i>	1.3135	0.786	0.9898	0.442	0.560
<i>Stipa pennata</i>	-3.59539	<b>0.018</b>	0.21154	<b>&lt;0.00001</b>	0.442
<i>Thymus vulgaris</i>	-3.0526	0.518	1.0597	0.687	0.603



# MOLECULAR ECOLOGY RESOURCES

**Table S3.** Correction factors for each species using one control species (three first columns, following Thomas et al. 2016 Mol. Ecol. Resources 16), and using all species mixtures (last column).

Species	Correction index using <i>Linum</i> as control species	Correction index using <i>Stipa</i> as control species	Correction index using <i>Helianthemum</i> as control species	Correction index based on all pairwise mixtures
<i>Bupleurum fruticosens</i>	2.4	0.192	1.026	0.817
<i>Helianthemum cinereum</i>	0.056	0.011	--	0.188
<i>Stipa pennata</i>	3.22	--	88.053	15.506
<i>Linum suffruticosum</i>	--	0.047	5.845	0.904
<i>Thymus vulgaris</i>	0.728	0.021	13.116	0.985

# MOLECULAR ECOLOGY RESOURCES

**Table S4.** Comparison of the percentage of biomass of each species in the multispecies mixtures (actual biomass) to the observed (uncorrected) and the corrected (after recalculation of the number of reads with correction indices) percentage of reads. Correction indices were computed using correction indices calculated with one control species (*Linum*, *Stipa* or *Helianthemum*, see Table S3), or using all species. The average error is calculated for the three replicates of each multispecies mixture as the absolute difference between the corrected and uncorrected biomass percentages for each species.

Species	<i>B. fruticescens</i>	<i>H. cinereum</i>	<i>L. suffruticosum</i>	<i>S. pennata</i>	<i>Th. vulgaris</i>	Average error
<b>Mixture 1. Actual Biomass (%)</b>	5.00	5.00	5.00	80.00	5.00	
Uncorrected biomass (%)	27.62	46.98	7.76	12.82	4.84	26.94
Corrected biomass ( <i>Linum</i> as control sp.)	54.15	2.27	6.42	34.32	2.84	20.23
Corrected biomass ( <i>Stipa</i> as control sp.)	33.75	3.12	1.71	60.86	0.56	11.50
Corrected biomass ( <i>Helianthemum</i> as control sp.)	3.44	4.77	3.31	82.55	5.94	1.40
Corrected biomass (using all spp.)	13.72	4.58	2.94	76.66	2.09	3.49
<b>Mixture 2. Actual Biomass (%)</b>	<b>22.50</b>	<b>22.50</b>	<b>22.50</b>	<b>10.00</b>	<b>22.50</b>	
Uncorrected biomass (%)	10.59	58.05	18.53	0.93	11.91	14.22
Corrected biomass ( <i>Linum</i> as control sp.)	44.14	6.16	30.90	4.92	13.88	12.02
Corrected biomass ( <i>Stipa</i> as control sp.)	43.23	15.00	17.98	18.67	5.12	11.76
Corrected biomass ( <i>Helianthemum</i> as control sp.)	2.75	16.36	26.09	19.19	35.61	10.35
Corrected biomass (using all spp.)	14.05	19.27	28.60	21.93	16.15	7.21
<b>Mixture 3. Actual Biomass (%)</b>	<b>20.00</b>	<b>20.00</b>	<b>20.00</b>	<b>20.00</b>	<b>20.00</b>	
Uncorrected biomass (%)	18.21	34.11	33.42	0.89	13.09	11.07
Corrected biomass ( <i>Linum</i> as control sp.)	45.90	2.16	37.67	3.08	11.18	17.43
Corrected biomass ( <i>Stipa</i> as control sp.)	49.19	6.14	26.36	13.21	5.10	14.22
Corrected biomass ( <i>Helianthemum</i> as control sp.)	4.19	7.48	38.56	17.17	32.59	12.46
Corrected biomass (using all spp.)	18.78	8.25	40.80	17.45	14.73	8.32
<b>Mixture 4. Actual Biomass (%)</b>	15.00	15.00	15.00	40.00	15.00	
Uncorrected biomass (%)	4.00	68.16	8.66	1.13	18.04	22.48
Corrected biomass ( <i>Linum</i> as control sp.)	23.98	11.04	22.41	9.74	32.82	13.69

# MOLECULAR ECOLOGY RESOURCES

Corrected biomass ( <i>Stipa</i> as control sp.)	21.43	23.38	11.22	32.84	11.14	5.92
Corrected biomass ( <i>Helianthemum</i> as control sp.)	0.90	16.46	12.11	22.88	47.65	13.64
Corrected biomass (using all spp.)	5.43	22.96	14.56	30.26	26.79	7.90

<b>Mixture 5. Actual Biomass (%)</b>	10.00	10.00	10.00	60.00	10.00	
Uncorrected biomass (%)	1.67	79.19	7.36	2.48	8.14	27.91
Corrected biomass ( <i>Linum</i> as control sp.)	13.40	19.16	23.21	25.13	19.10	13.95
Corrected biomass ( <i>Stipa</i> as control sp.)	7.73	26.70	7.49	53.99	4.10	6.68
Corrected biomass ( <i>Helianthemum</i> as control sp.)	0.38	22.49	8.91	45.17	23.05	10.21
Corrected biomass (using all spp.)	2.00	26.88	9.62	51.01	10.50	6.95

<b>Mixture 6. Actual Biomass (%)</b>	22.50	10.00	22.50	22.50	22.50	
Uncorrected biomass (%)	9.59	51.86	28.36	2.46	7.45	19.15
Corrected biomass ( <i>Linum</i> as control sp.)	34.00	5.27	41.04	11.97	7.72	12.01
Corrected biomass ( <i>Stipa</i> as control sp.)	28.50	10.52	20.19	38.35	2.43	8.95
Corrected biomass ( <i>Helianthemum</i> as control sp.)	1.79	11.09	29.71	40.04	17.38	10.33
Corrected biomass (using all spp.)	8.55	12.25	30.02	41.95	7.24	11.69

<b>Mixture 7. Actual Biomass (%)</b>	15.00	40.00	15.00	15.00	15.00	
Uncorrected biomass (%)	6.38	86.18	2.44	0.18	4.81	18.47
Corrected biomass ( <i>Linum</i> as control sp.)	41.56	23.66	12.26	3.20	19.33	12.35
Corrected biomass ( <i>Stipa</i> as control sp.)	34.95	45.37	5.93	8.22	5.53	10.13
Corrected biomass ( <i>Helianthemum</i> as control sp.)	4.48	47.75	7.31	8.93	31.52	9.71
Corrected biomass (using all spp.)	15.87	52.56	8.03	8.99	14.55	5.37

<b>Mixture 8. Actual Biomass (%)</b>	10.00	60.00	10.00	10.00	10.00	
Uncorrected biomass (%)	4.15	90.93	2.53	0.32	2.07	12.37
Corrected biomass ( <i>Linum</i> as control sp.)	33.85	34.69	17.44	4.40	9.63	12.51
Corrected biomass ( <i>Stipa</i> as control sp.)	24.42	55.27	6.42	11.63	2.26	6.42
Corrected biomass ( <i>Helianthemum</i> as control sp.)	2.31	56.92	9.06	15.16	16.55	4.68
Corrected biomass (using all spp.)	9.39	61.01	8.84	14.28	6.47	2.12

# MOLECULAR ECOLOGY RESOURCES

---

<b>Mixture 9. Actual Biomass (%)</b>	5.00	80.00	5.00	5.00	5.00	
Uncorrected biomass (%)	0.91	96.85	1.26	0.05	0.94	6.74
Corrected biomass ( <i>Linum</i> as control sp.)	20.27	58.97	12.65	1.34	6.78	9.88
Corrected biomass ( <i>Stipa</i> as control sp.)	11.67	79.79	4.13	2.99	1.43	2.67
Corrected biomass ( <i>Helianthemum</i> as control sp.)	0.73	80.19	5.92	3.31	9.85	2.38
Corrected biomass (using all spp.)	3.27	84.08	5.61	3.17	3.86	1.88

---

# MOLECULAR ECOLOGY RESOURCES

**Figure S2.** Average error (average absolute difference between the actual percentage of biomass and the corrected/uncorrected percentages of reads) for each multispecies mixture (averaged across species). Corrections were performed based on different correction indices (see Table S3). Species-specific corrected biomass can be found on Table S4. The composition and relative abundance of each multispecies mixture (code numbers in the X axis match those on Table S4).

