



# Widespread sampling biases in herbaria revealed from large-scale digitization

The Harvard community has made this article openly available. [Please share](#) how this access benefits you. Your story matters

Citation	Daru, Barnabas H., Daniel S. Park, Richard B. Primack, Charles G. Willis, David S. Barrington, Timothy J. S. Whitfeld, Tristram G. Seidler, et al. 2017. "Widespread Sampling Biases in Herbaria Revealed from Large-Scale Digitization." <i>New Phytologist</i> (October 30). doi:10.1111/nph.14855.
Published Version	doi:10.1111/nph.14855
Citable link	<a href="http://nrs.harvard.edu/urn-3:HUL.InstRepos:34389682">http://nrs.harvard.edu/urn-3:HUL.InstRepos:34389682</a>
Terms of Use	This article was downloaded from Harvard University's DASH repository, and is made available under the terms and conditions applicable to Open Access Policy Articles, as set forth at <a href="http://nrs.harvard.edu/urn-3:HUL.InstRepos:dash.current.terms-of-use#OAP">http://nrs.harvard.edu/urn-3:HUL.InstRepos:dash.current.terms-of-use#OAP</a>

1 **Widespread sampling biases in herbaria revealed from large-scale digitization**

2

3

4

5 Barnabas H. Daru<sup>1, \*</sup>, Daniel S. Park<sup>1, \*</sup>, Richard B. Primack<sup>2</sup>, Charles G. Willis<sup>1</sup>, David S.  
6 Barrington<sup>3</sup>, Timothy J. S. Whitfeld<sup>4</sup>, Tristram G. Seidler<sup>5</sup>, Patrick W. Sweeney<sup>6</sup>, David  
7 R. Foster<sup>7</sup>, Aaron M. Ellison<sup>7,8</sup> and Charles C. Davis<sup>1</sup>

8

9 <sup>1</sup>Department of Organismic and Evolutionary Biology, Harvard University Herbaria,  
10 Harvard University, Cambridge, MA 02138, USA; <sup>2</sup>Biology Department, Boston  
11 University, Boston, MA 02215, USA; <sup>3</sup>Pringle Herbarium, Plant Biology Department,  
12 University of Vermont, Torrey Hall, 27 Colchester Ave, Burlington, VT 05405, USA;  
13 <sup>4</sup>Brown University Herbarium, Department of Ecology and Evolutionary Biology, Brown  
14 University, 34 Olive Street, Box G-B225, Providence, Rhode Island 02912 USA;  
15 <sup>5</sup>Biology Department, University of Massachusetts, 611 North Pleasant Street, Amherst,  
16 MA 01003, USA; <sup>6</sup>Division of Botany, Peabody Museum of Natural History, Yale  
17 University, New Haven, CT, USA; <sup>7</sup>Harvard Forest, Harvard University, 324 North Main  
18 Street, Petersham, Massachusetts, 01366 USA; <sup>8</sup>Tropical Forests & People Research  
19 Centre, University of the Sunshine Coast, Maroochydore, Queensland 4558, Australia.

20

21 \*These authors contributed equally to the study

22

23 <sup>1</sup>To whom correspondence should be addressed. Email: [barnabas\\_daru@fas.harvard.edu](mailto:barnabas_daru@fas.harvard.edu)

24 Tel: +1 857 218 0117

25 Twitter: @Barnabas\_Daru

26

27 **Short Title:** Sampling bias in herbarium specimens

28 **Manuscript information:** 5872 words (Introduction = 731 words, Materials and  
29 Methods = 1651 words, Results = 1067 words, Discussion = 2423 words | 8 figures (8  
30 color figures) | 3 Tables | 1 supporting information

31

32 **SUMMARY**

- 33 1. Non-random collecting practices may bias conclusions drawn from analyses of  
34 herbarium records. Recent efforts to fully digitize and mobilize regional floras online  
35 offer a timely opportunity to assess commonalities and differences in herbarium  
36 sampling biases.
- 37 2. We determined spatial, temporal, trait, phylogenetic, and collector biases in ~5  
38 million herbarium records, representing three of the most complete digitized floras of  
39 the world: Australia (AU), South Africa (SA), and New England, USA (NE).
- 40 3. We identified numerous shared and unique biases among these regions. Shared biases  
41 included specimens i) collected close to roads and herbaria; ii) collected more  
42 frequently during biological spring and summer; iii) of threatened species collected  
43 less frequently; and iv) of close relatives collected in similar numbers. Regional  
44 differences included i) over-representation of graminoids in SA and AU and of  
45 annuals in AU; and ii) peak collection during the 1910s in NE, 1980s in SA, and  
46 1990s in AU. Finally, in all regions, a disproportionately large percentage of  
47 specimens were collected by very few individuals. We hypothesize that these mega-  
48 collectors, and along with their associated preferences and idiosyncrasies, shaped  
49 patterns of collection bias via ‘founder effects’.
- 50 4. Studies using herbarium collections should account for sampling biases, and future  
51 collecting efforts should avoid compounding these biases to the extent possible.

52 **Keywords:** Herbarium, collector bias, geographic bias, regional flora, sampling bias,  
53 temporal bias, trait bias

54

55 **INTRODUCTION**

56 Herbaria contain a wealth of information about the ecological and evolutionary history of  
57 living and extinct species (Funk, 2003). Despite the continuous decline in plant collecting  
58 and declining support for herbaria (Dalton, 2003; Prather *et al.*, 2004a, b), there has been  
59 a recent surge of studies leveraging herbarium collections for diverse research projects  
60 not focused on systematics (Pyke & Ehrlich, 2010; Lees *et al.*, 2011; Feeley, 2012;  
61 Lavoie, 2013; Hart *et al.*, 2014). These studies include plant demography, current and  
62 future species distributions, and temporal changes in phenology and morphology (*e.g.*,  
63 Miller-Rushing *et al.*, 2006; Newbold, 2010; Pyke & Ehrlich, 2010; Lavoie, 2013; Staats  
64 *et al.*, 2013; Davis *et al.*, 2015; Willis *et al.*, 2017a,b).

65 Ideally, herbarium collections used for these studies would include statistically  
66 unbiased samples of plant diversity across space and time. However, as the majority of  
67 specimens were collected for qualitative taxonomic and/or systematic inquiries, they  
68 were usually collected non-randomly and sampling designs were rarely quantified (Wolf  
69 *et al.*, 2011; Schmidt-Lebuhn *et al.*, 2013). Because non-random samples may be  
70 statistically biased, analyzing them without accounting for biases might lead to spurious  
71 results (Syfert *et al.*, 2013).

72 Sampling biases fall into several broad categories. Taxonomic or phylogenetic  
73 bias is the unbalanced sampling of certain taxa or clades over others, typically resulting  
74 from the scientific interests of a collector or the attractiveness of plants (Hortal *et al.*,  
75 2007). Geographic bias occurs when specimens are collected more frequently in one  
76 place than another, often because of differential accessibility (Hijmans *et al.*, 2000).  
77 Temporal bias occurs when collection activity is favored in certain years or parts of the  
78 year (Cotterill *et al.*, 1994; Funk & Morin, 2000; Norris *et al.*, 2001). Meyer *et al.* (2016)  
79 evaluated worldwide terrestrial plant occurrence data using 120 million records from the  
80 Global Biodiversity Information Facility (GBIF; Edwards *et al.*, 2000). Their analyses  
81 revealed large taxonomic gaps in global plant occurrence data (< 25% of species of land  
82 plants were sampled); extensive spatial gaps across regions that harbor high  
83 concentrations of plant diversity, especially in Asia, Central Africa, and Amazonia; and  
84 strong temporal discontinuities in occurrence records across decades, all of which can  
85 hamper inferences about the effects on plants of recent and future environmental change.

86           Although Meyer *et al.*'s (2016) study represents the most comprehensive effort to  
87 assess biases in plant collections at a global scale to date, the vast majority of herbarium  
88 collections have not been digitized, and of those that have, many are unavailable, in  
89 whole or in part, on GBIF. Thus, Meyer *et al.*'s (2016) assessment of biases may itself be  
90 biased, or may inaccurately reflect biases in more complete, regional botanical  
91 collections that have been more fully mobilized. Furthermore, over two-thirds of the plant  
92 records in GBIF are not tied to physical specimens, and thus cannot be easily validated by  
93 others (Cotterill, 1995). For these reasons we suspect that an analysis of finer-grained  
94 collection data, focused on specific regions that have been predominantly digitized and  
95 validated, may reveal clearer patterns of sampling biases between regions than the global  
96 trends identified by Meyer *et al.* (2016) (*cf.* Hijmans *et al.*, 2000 for Bolivian potatoes).

97           Expanding upon Meyer *et al.*'s work, we explored spatial, temporal, and  
98 taxonomic/phylogenetic sampling biases in collections from three of the most extensively  
99 collected, digitized, and mobilized regional floras in the world: South Africa (SA),  
100 Australia (AU), and the New England (NE) region of the United States. The SA flora is a  
101 compilation of digitized herbarium specimens from all major herbaria across the country  
102 available in a single online portal (South African National Biodiversity Institute [SANBI],  
103 2016; le Roux *et al.*, 2017). The Australian Virtual Herbarium (AVH, 2016) is the main  
104 database for AU. It contains digitized herbarium specimens from all the major herbaria in  
105 AU. The Consortium of the Northeast Herbaria database contains digitized specimens  
106 from 15 participating herbaria in the NE region of the United States (Schorn *et al.*, 2016).  
107 We also examined trait bias – sampling bias due to intrinsic life-history characteristics,  
108 including life cycle (annual *vs.* perennial), plant height, growth form (woody *vs.*  
109 herbaceous), and species conservation status. Finally, we examined the contributions of  
110 individual collectors to each flora. We identified biases in all five of these categories  
111 within each of these regional floras. Our results revealed both commonalities and  
112 differences in regional collection biases and identified new sampling foci as collections  
113 grow in the future.

114

## 115 **MATERIAL AND METHODS**

### 116 **Sources and description of data**

117 We obtained 12,488,200 herbarium specimen records of vascular plants from AU  
118 (Australia Virtual Herbarium [AVH], 2016); 2,049,905 herbarium specimen records from  
119 SA including Lesotho and Swaziland (South African National Biodiversity Institute  
120 [SANBI], 2016); and 879,388 herbarium specimen records from the NE (USA) flora  
121 (Consortium of Northeastern Herbaria [CNH], 2016). The records were cleaned in two  
122 steps (Fig. S1). First, we standardized the taxonomy of all species using the Taxonomic  
123 Name Resolution Service v.4.0 (Boyle *et al.*, 2013). This online tool corrects and  
124 standardizes plant names against reference taxonomies, such as Missouri Botanical  
125 Garden's Tropicos (<http://tropicos.org/>) database or the PlantList (<http://theplantlist.org/>).  
126 Second, we removed specimens that were duplicates from the same collection locality  
127 and date; specimens with clearly erroneous locations (*i.e.*, in oceans); specimens with  
128 zero coordinates and occurrences that fell outside the boundaries of our study; specimens  
129 missing exact collection date or georeferenced location data; and field observation  
130 records not tied to a physical specimen. Following this data cleaning, we retained 32% of  
131 the initial specimens for further analysis: 24% of the AU records (31,966 taxa; 2,958,195  
132 records); 49% (20,824 taxa, 1,008,206 records) from SA; and 75% (3719 taxa, 661,370  
133 records) from NE.

134

## 135 **Analyses**

### 136 ***Spatial biases***

137 First, we evaluated the density of sampling localities across the focal regions using  
138 Delaunay triangulation polygons, which measure the land area covered by each sampling  
139 locale (Fortune, 1992). Larger triangles indicate sparser collecting effort, whereas smaller  
140 triangles indicate more concentrated effort. Second, we examined infrastructure bias by  
141 calculating the minimum distance of each collection locality to the nearest major road  
142 (GADM, 2015) and herbarium (following Thiers, 2016). Our dataset of roads derives  
143 from the publicly available Digital Chart of the World (<http://maproom.psu.edu/dcw/>),  
144 which was compiled by the US Defense Mapping Agency from 1:1,000,000 scale paper  
145 maps (ESRI, 1992). All roads appearing at this scale were included in our analyses.  
146 Although this dataset includes only larger roads and has not been updated since 1992, it  
147 likely represents the most comprehensive digital record of roads around the world. We

148 then compared these distances to those generated by a null model (1000 iterations) in  
149 which the same number of sample points was randomly (Poisson) distributed across each  
150 geographic region. Third, we mapped geographic biases in sampling density, defined as  
151 areas of excessive (hotspots) or insufficient (coldspots) collection (Hijmans *et al.*, 2000).  
152 Hotspots and coldspots were determined at a spatial grain of  $0.25^\circ \times 0.25^\circ$  based on the  
153 number of specimens per grid cell, and identified using the 2.5% threshold (Ceballos &  
154 Ehrlich, 2006; Orme *et al.*, 2005; Daru *et al.*, 2015), based, respectively, on the 97.5<sup>th</sup> and  
155 2.5<sup>th</sup> percentile values in the number of specimens collected per grid cell. Spatial distance  
156 calculations were computed with the functions *dist2Line* and *spDists* in the R packages *sp*  
157 (Bivand *et al.*, 2013) and *geosphere* (Hijmans, 2015), respectively. In our final predictive  
158 model of sampling density, we also included human population density (CIESIN, 2016),  
159 sampling localities, infrastructure (distance to herbaria and roads), number of specimens  
160 collected, and elevation

161

### 162 ***Temporal bias***

163 For each regional flora, we explored bias at several temporal scales. Collection dates  
164 ranged from 20 May 1664 to 9 January 2016 (AU), 15 November 1656 to 6 June 2016  
165 (SA), and 28 July 1687 to 4 May 2016 (NE). We hypothesized that collectors tended to  
166 avoid fieldwork during unfavorable conditions (*e.g.*, winter, wartime) or certain days of  
167 the week (*e.g.*, weekdays for non-professional botanists). To test for temporal bias, we  
168 first re-coded collection dates as days of the week (Sunday = 1, Monday = 2, *etc.*), and  
169 day of the year (DOY; where January 1 = 1 DOY and December 31 = 365 DOY, *etc.*).  
170 We then used a Rayleigh test of directional statistics in the R package *circular*  
171 (Agostinelli & Lund, 2013) to test whether each of these collection dates were randomly  
172 distributed against all dates spanning the entire duration of plant collection. If  $P < \alpha =$   
173 0.05, we rejected the null hypothesis of temporal uniformity at scales of weeks, days of  
174 the year, or decades.

175

### 176 ***Trait bias***

177 We used customized R scripts to harvest information on growth duration (annual *vs.*  
178 perennial), growth form (woody *vs.* herbaceous), and height for each species from online

179 regional databases (all accessed in June 2016), including: New South Wales Flora Online  
180 (<http://plantnet.rbgsyd.nsw.gov.au>); JSTOR Global Plants (<https://plants.jstor.org>); Atlas  
181 of Living Australia (<http://bie.ala.org.au>); Plants of Southwestern Australia  
182 (<http://keys.lucidcentral.org>); the African Plant Database (<http://ville-ge.ch>); Plants of  
183 Southern Africa (<http://plantzafrica.com>); Plant Resources of Tropical Africa  
184 (<http://www.prota4u.org>); Flora of North America (<http://efloras.org>); and the USDA  
185 Plants Database (<http://plants.usda.gov>). We then manually checked these data for  
186 inconsistencies in terminologies for defining certain traits. For example, ‘vines’ vs.  
187 ‘lianas’ for climbers, ‘forbs’ vs. ‘herbs’ for herbaceous life forms, ‘biennial’ for perennial  
188 growth duration. Extinction risk assessments for each species were obtained from the  
189 IUCN Red List database ([www.iucnredlist.org](http://www.iucnredlist.org), accessed August 2016), which uses the  
190 following categories: Data Deficient (DD), Least Concern (LC), Lower  
191 Risk/Conservation Dependent (LR/CD), Near Threatened (NT), Vulnerable (VU),  
192 Endangered (EN), Critically Endangered (CR), and Extinct (EX). We grouped these  
193 narrow categories into two broader threat categories, threatened (EX+CR+EN+VU) or  
194 not threatened (LR/CD+NT+LC), following Yessoufou *et al.* (2012).

195 Trait bias was evaluated using a chi-squared test to contrast the number of  
196 observed specimens collected per species with the abundance of a species if specimen  
197 collection was equal across all species for each trait category. Because of dramatically  
198 unequal sampling effort in some species – *e.g.*, *Senna artemisioides* with 10,167  
199 specimens vs. *Eucalyptus cordieri* with only one – and the low coverage of taxa with  
200 available trait data, we randomly sampled 50 specimens from each available species with  
201 trait data using 1000 randomizations. Species with less than 50 specimens were excluded  
202 from this analysis.

203

#### 204 ***Phylogenetic bias***

205 We assessed phylogenetic signal in collection frequency as a measure of phylogenetic  
206 bias using two different tests following Wolkovich *et al.* (2013). A strong phylogenetic  
207 signal – closely related species sharing similar collection frequency – would suggest  
208 phylogenetic bias in collections. We first assembled a phylogeny using Phylomatic  
209 (Webb & Donoghue, 2005), enforcing a topology that assumed the APG III (APG III,



210 2009) backbone (tree R20120829). This phylogeny included all species in our analysis,  
211 but provided only an approximate degree of relatedness based on taxonomic hierarchy at  
212 family level; many relationships, especially within genera, were unresolved. This is  
213 problematic because recent theoretical and empirical studies have shown that a lack of  
214 resolution in a community phylogeny may mask significant patterns by reducing  
215 statistical power (Schaefer *et al.*, 2011; Daru *et al.*, 2017) or suggest significant  
216 phylogenetic patterns that are not supported by more completely resolved phylogenies  
217 (Davies *et al.*, 2012).

218 To alleviate these concerns, we also tested for phylogenetic bias by including only  
219 those species sampled in the dated molecular phylogeny inferred from seven genes for  
220 32,223 plant species (Zanne *et al.*, 2014). Although this phylogeny has been criticized  
221 (Edwards *et al.*, 2015), it nonetheless represents the single largest phylogeny to date for  
222 flowering plants. The taxon sampling for testing phylogenetic bias included 5814 species  
223 from AU, 3568 from SA, and 4269 from NE.

224 We estimated phylogenetic signal using three common metrics: Abouheif's  $C_{\text{mean}}$   
225 (Abouheif, 1999), Blomberg's  $K$  (Blomberg *et al.*, 2003), and Pagel's  $\lambda$  (Pagel,  
226 1999). Significance was assessed by comparing observed values to a null distribution  
227 created by shuffling the trait values across the tips of the phylogeny 1000 times. Pagel's  $\lambda$   
228 uses a maximum-likelihood method with branch-length transformation to estimate the  
229 best-fit of a trait against a Brownian model. Values of Pagel's  $\lambda$  range from 0 (no  
230 phylogenetic signal) to 1 (strong phylogenetic signal). Both Blomberg's  $K$  (a significant  
231 phylogenetic signal is indicated by a  $K$  value  $> 1$ ) and Pagel's  $\lambda$  were calculated using the  
232 R package *phytools* (Revell, 2012). Abouheif's  $C_{\text{mean}}$  was calculated using *adephylo*  
233 (Jombart & Dray, 2008). We tested the sensitivity of our analysis by exploring  
234 phylogenetic signal in collecting effort across nine well-sampled clades as represented in  
235 NE: Asteraceae, Brassicaceae, Cyperaceae, Ericaceae, Fabaceae, Lamiaceae, Poaceae,  
236 Ranunculaceae, and Rosaceae.

237 In addition to phylogenetic signal, we also used phylogenetic generalized least  
238 squares regressions (PGLS) in the R package *caper* (Orme *et al.*, 2012) to model  
239 collecting effort per species in each region as a function of species evolutionary ages,  
240 evolutionary distinctiveness (ED), and 'evolutionary distinctiveness and global

241 endangerment' (EDGE; Isaac *et al.*, 2007). Species ages were measured as the length of  
242 terminal branches (BL) linking species on a phylogenetic tree. ED measures the degree of  
243 phylogenetic isolation of a species, whereas the EDGE metric was determined by  
244 calculating the ED score of each species (Isaac *et al.*, 2007) and combining it with global  
245 endangerment (GE) from IUCN conservation categories:  $EDGE = \ln(1 + ED) + GE \times$   
246  $\ln(2)$ , where GE represents expected probability of species extinction over a 100-year  
247 period (Redding & Mooers, 2006) categorized as follows: least concern = 0.001, near  
248 Threatened and Conservation Dependent = 0.01, Vulnerable = 0.1, Endangered = 0.67,  
249 and Critically Endangered = 0.999.

250 Last, we examined the phylogenetic structure of collecting efforts across decades  
251 to test for patterns of phylogenetic overdispersion and clustering through time. Temporal  
252 phylogenetic structure by decade (*i.e.*, 1901-1910, 1911-1920, *etc.*) was evaluated using  
253 the net relatedness index (NRI) and nearest taxon index (NTI; Webb *et al.*, 2002, 2008).  
254 NRI describes a tree-wide pattern of phylogenetic dispersion, whereas NTI evaluates  
255 phylogenetic structure towards the tips of the phylogeny. Negative values of NRI or NTI  
256 indicate phylogenetic overdispersion whereas positive values indicate phylogenetic  
257 clustering.

258

### 259 ***Collector bias***

260 We determined collector bias by tabulating the number of specimens amassed by each  
261 collector in all three floras. We then examined Pearson's product-moment correlation  
262 between the numbers of specimens collected per collector with the number of species  
263 collected per collector.

264

### 265 **Computation and availability of data and code**

266 All statistical analyses were conducted using the Research Computing Clusters of  
267 Harvard University (<https://rc.fas.harvard.edu/>). Data files and custom R scripts are  
268 available from the Harvard Forest Data Archive, dataset HF296  
269 (<http://harvardforest.fas.harvard.edu/data-archive>).

270

## 271 **RESULTS**

272 ***Spatial bias***

273 High sampling density was observed in southeast and southwest AU, the Cape region of  
274 SA, and two of the six NE states (Connecticut and Massachusetts) relative to other parts  
275 of those regions (Fig. 1a-c). When we weighted each sampling locale by the number of  
276 specimens, we found a mismatch between hotspots (top 2.5% quantiles) and coldspots  
277 (lowest 2.5% quantiles) of sampling intensity (Fig. 1d-f), suggesting hotspots and  
278 coldspots were not randomly distributed. Hotspots of collecting tend to cluster around  
279 coasts in AU and SA, whereas coldspots were abundant in interior areas. In NE, hotspots  
280 were concentrated in the south and coldspots occurred in the north.

281 Herbarium specimens tended to be collected closer than expected to roads and  
282 herbaria ( $p < 0.01$ ; Fig. 2a, b). More than 50% of herbarium specimens were collected  
283 within 2 km of roadsides in all three floras ( $p < 0.01$ ; Fig. 2a). Moreover, distance to  
284 herbaria explained 45% of the variance in collecting effort in AU, 29% in SA and 12.3%  
285 in NE, with a higher density of specimens closer to herbaria (Table 1). Despite substantial  
286 gradients in altitudes in each region (-15 – 2022 m a.s.l. in AU; 1 – 3254 m a.s.l. in SA;  
287 and -3 – 1485 m a.s.l. in NE), most specimens were collected below 500 m a.s.l. in AU  
288 and NE (81%, 44%, and 93% of specimens in AU, SA, and NE, respectively; Fig. 2c).  
289 We also found a negative correlation between collecting effort and altitude in Australia  
290 and South Africa, suggesting a tendency for specimens to be collected at lower elevations.  
291 However, the opposite was true for NE, where more specimens tended to be collected at  
292 higher elevations than expected by chance (Table 1).

293

294 ***Temporal bias***

295 There were historical biases in collection efforts in the three floras: low sampling until  
296 1880 in AU and SA, and a burst of collections in NE in the early 20<sup>th</sup> century (Fig. 3).  
297 Conversely, there was a dramatic increase in botanical collection in SA and AU after  
298 World War II, peaking in the 1980s and 1990s, respectively (Fig. 3). This peaking  
299 occurred ~100 years after peak collection activity in NE. Seasonally, specimen  
300 collections were biased toward spring and summer for all three floras, with peak  
301 collections ranging from September to December in AU and SA (Rayleigh  $Z = 0.189$  and  
302  $Z = 0.251$  respectively, both  $p < 0.001$ ), and May to September in NE (Rayleigh  $Z =$

303 0.718,  $p < 0.001$ ; Fig. 4a). There was a significant trend towards collection on weekends  
304 (Saturdays and Sundays) in NE (Rayleigh test  $Z = 1.0$ ,  $p < 0.001$ ) and midweek in SA  
305 and AU (Rayleigh test  $Z = 0.105$  and  $Z = 1.0$ , respectively; both  $p < 0.001$ ; Fig. 4a).

306

### 307 ***Trait bias***

308 Perennials were more frequently collected than annuals in terms of specimens per species  
309 in SA and NE; the opposite was true for AU where there was a greater tendency for  
310 annuals to be collected (Fig. 5a). Similarly, graminoid specimens per species were over-  
311 represented relative to other habits in AU and SA, whereas herbs and trees were over-  
312 represented in NE (Fig. 5b). Relatively short plants were more frequently represented  
313 than taller plants in all three floras: 79.3%, 89.3% and 84.9% of the plants collected in  
314 AU, SA and NE, respectively were less than 5 m in height (Fig. 5c).

315 Threatened species were collected significantly less often than non-threatened  
316 plants across all three floras (all  $p < 0.001$ ; Fig. 5d).

317

### 318 ***Phylogenetic bias***

319 It is possible that closely related species were collected similarly (either under-collected  
320 or over-collected) than expected by chance. We tested this assumption phylogenetically  
321 and found a significant, but weak phylogenetic signal in the abundance of specimens per  
322 species across all three floras (Table 2). Specifically, closely related species tended to  
323 have a more similar number of specimens than expected (Table 2; Fig. 6). This  
324 phylogenetic bias was strongest in SA (Abuohelf's  $C_{\text{mean}} = 0.15$  and  $\lambda = 0.32$ ; both  $p <$   
325  $0.01$ , but  $K = 0.0013$  [NS]). For instance, in SA, collections from the genus *Protea*  
326 averaged 115 specimens per species whereas only two specimens were collected for  
327 species in the genus *Rytigynia* on average. Most *Agoseris* in NE were represented by  $< 10$   
328 specimens per species, whereas many fern genera were represented by high specimen  
329 numbers (e.g., *Onoclea* with 845 specimens/species). Australian collections showed the  
330 weakest phylogenetic bias (Abuohelf's  $C_{\text{mean}} = 0.12$  and  $\lambda = 0.18$ , both  $p < 0.01$ , but  $K =$   
331  $0.00085$  [NS]; Fig. 6). Phylogenetic signal varied at the family level as well in NE, with  
332 Asteraceae showing the strongest collection bias (Fig. 7), followed by Cyperaceae,

333 Poaceae, and Rosaceae (Table S1). These families were represented by much higher  
334 collection numbers than for other families.

335 EDGE was a significant predictor of collecting efforts in all three floras ( $p <$   
336 0.001), with variance ranging from 1.89% (NE) and 3.75% (AU), to 8.89% in SA. In  
337 general, EDGE species (distinctive, endangered species) were generally under-collected  
338 in terms of specimens per species (Table 3).

339 Lastly, floristic collecting showed a general trend of phylogenetic clustering  
340 within decades for all three floras. The collection of different clades of plants was not  
341 evenly distributed across time. NTI was significantly positive in each flora, indicating  
342 that clustering occurred near the tips of the phylogeny (Fig. 3). We only observed  
343 significant phylogenetic clustering at the deeper nodes of the phylogeny, as indicated by  
344 NRI, in SA (Fig. 3); deeper phylogenetic clustering was weak in NE and AU (Fig. 3).

345

#### 346 *Collector bias*

347 The number of specimens per collector was highly skewed (Fig. 8). In AU, more than 50%  
348 of the examined specimens were amassed by only 2% of the collectors, including A.C.  
349 Beaglehole (46,728 specimens), B. Hyland (32,019 specimens), and P.I. Forster (30,280  
350 specimens; Fig. 8a). In SA, more than 50% of the specimens were amassed by 9.5% of  
351 collectors, including J.P.H Acocks (19,344 specimens), E.E. Esterhuysen (15,566  
352 specimens), and E.E. Galpin (14,146 specimens; Fig. 8b). In NE, 50% of the specimens  
353 were contributed by 3.2% of the collectors, including L.J. Mehrhoff (19,149 specimens),  
354 M.L. Fernald (14,368 specimens), and A.S. Pease (12,238 specimens; Fig. 8c). The  
355 number of specimens amassed by these collectors was very strongly positively correlated  
356 with the number of species they collected, suggesting that these collectors were mainly  
357 doing general collecting rather than focusing on a particular group of plants ( $r = 0.85$  in  
358 AU, 0.95 in SA and 0.84 in NE; all  $p < 0.01$ ; Fig. S2).

359

#### 360 **DISCUSSION**

361 Historically, the primary function of herbaria has been to serve as a resource for botanists  
362 carrying out taxonomic and systematic research, allowing users to construct  
363 classifications of plants, verify identifications, determine the ranges and morphological

364 characteristics of species, and develop local and regional floras (Greve *et al.*, 2016). Over  
365 time, new uses for specimens have arisen, and now more than ever, they are being used in  
366 ways that collectors rarely imagined (Pyke & Ehrlich, 2010; Lavoie, 2013; Willis *et al.*,  
367 2017a,b; Nualart *et al.*, 2017; Rudin *et al.*, 2017). Accordingly, attempts to assess and  
368 categorize biases inherent in these collections have been made (Rich & Woodruff, 1992;  
369 Geri *et al.*, 2013; Schmidt-Lebuhn *et al.*, 2013; Meyer *et al.*, 2016; Stropp *et al.*, 2016).  
370 Among these, the most comprehensive investigation is by Meyer *et al.* (2016), who  
371 proposed an important conceptual framework for analyzing gaps and biases along  
372 taxonomic, geographical, and temporal dimensions. Although Meyer *et al.* (2016)  
373 focused more on observational records than herbarium collections, they uncovered  
374 numerous biases in ‘digitally accessible information’ regarding plants and provided an  
375 important baseline for evaluating and improving global floristic coverage in collection  
376 data. However, collection biases in large geographic areas are difficult to categorize  
377 because the collections data are not yet digitized, and this may skew the global patterns  
378 of bias noted by Meyer *et al.* (2016). By focusing on three of the most well-collected and  
379 digitized floras in the world, we reduced effects of missing or unavailable data, and most  
380 importantly, could evaluate commonalities and differences in patterns of bias among  
381 regional collections.

382

### 383 ***Spatial bias***

384 Our data confirmed the tendencies for botanists to collect along roadsides (*e.g.*, Funk &  
385 Richardson, 2002), near herbaria (*e.g.*, Hijmans *et al.*, 2000; Moerman & Estabrook,  
386 2006; Pautasso & McKinney, 2007), in more accessible areas (Rich & Woodruff, 1992),  
387 and at lower elevations. Before automobiles became common in the 1920s, botanists  
388 often walked or rode domesticated animals to collection sites (Botkin, 1968; Belasco,  
389 1979). As our modern infrastructure developed (*e.g.*, roads, railroads, and cities that  
390 contain herbaria) and access to automobiles increased, spatial biases associated with  
391 infrastructure likely increased as botanists tended to travel and collect using automobiles  
392 (Everill *et al.*, 2014). Because roads are known alter local environmental conditions and  
393 facilitate biological invasions (*e.g.*, Forman & Alexander, 1998; Hui *et al.*, 2003; Griffith

394 *et al.*, 2010; Li *et al.*, 2014) and botanists and herbaria predominate in cities, specimens  
395 collected in proximity to either are unlikely to represent a random sample across species  
396 distributions. Specifically, species collected along roadsides are likely to be over-  
397 represented by species that thrive with disturbance, and under-represented by forest  
398 interior and wetland species that are harmed by disturbance (Gutzwiller & Flather, 2011;  
399 Rivers-Moore & Cowden, 2012). As the road network continues to expand and as people  
400 become evermore concentrated in cities, this bias toward collecting near roads might  
401 become stronger in coming decades.

402         Collection bias towards lower elevations (< 500 m) was striking in SA and AU,  
403 despite extensive collection efforts in adjacent hyper-diverse hotspots such as the  
404 mountains in the Cape Fold Belt (SA), and Mount Lesueur-Eneabba (Western AU). This  
405 is likely due to the presence of the arid and relatively species-poor Great Karoo Plateau  
406 (SA), Great Sandy Desert (AU), which each encompass over a third of the respective  
407 study sites, but account for only a small proportion of the biodiversity of each region. As  
408 a result, the low-elevation collection bias in the floras may reflect actual species  
409 abundance. In NE, the trend toward collecting at higher elevation might be due to the  
410 strong tendency for botanists to visit the White Mountains and Mount Katadhin to collect  
411 alpine species.

412         Although we realize that patterns of species richness may not be randomly  
413 distributed across the landscape, accounting for underlying patterns of richness or  
414 abundance is difficult because our knowledge of such patterns often derive from (and are  
415 thus not independent from) these same (biased) collections. By comparing locations of  
416 samples (collections) against a Poisson set of points and specimens per species, and not  
417 total collection numbers, we tested only for the non-random distribution of collection  
418 locations on a landscape. And indeed, we found that the collection locations were not  
419 spatially distributed randomly (Poisson) on the landscape. It is also possible that  
420 georeferencing might have introduced additional bias in some specimens. While  
421 ascertaining the degree of accuracy of georeferenced records might be challenging  
422 because such information is often unavailable, our cumulative curves are likely less  
423 affected.

424

425 ***Temporal and seasonal bias***

426 Collections in AU and SA have increased through time until a few decades ago, but those  
427 in NE peaked much earlier in the early 1900s. These differences between regional  
428 collection activities may parallel broader societal factors influencing plant collection. In  
429 NE, for example, the establishment of the New England Botanical Club during the 1890s  
430 (NEBC, 1899) preceded a surge and peak in collecting activity associated with prolific  
431 botanical expeditions of the region coinciding with the ‘Golden Age’ of plant collecting  
432 in Europe and North America (Whittle, 1970; Musgrave *et al.*, 1999). In SA, collection  
433 efforts began much later, peaking during the Apartheid Era (1948–1994), and declined  
434 thereafter under the New Democratic Rule, concomitant with the general economic  
435 decline of the country and concern for public safety (Ferreira & Harmse, 2000; Lemanski,  
436 2004). In AU, the mass immigration of Europeans in 1948 after World War II included  
437 numerous highly skilled professionals (Price, 1998; Leuner, 2007) and coincided with an  
438 enormous increase in botanical collecting. Botanical collecting may have declined more  
439 recently owing to legislation in AU and SA to regulate collections activities, especially  
440 those designed to protect rare and endangered species.

441 Collecting efforts within a season revealed common patterns of bias: specimens in  
442 all three regions were collected overwhelmingly in biological spring and summer.  
443 Sampling during these time periods likely reflects efforts to collect plants in good  
444 flowering and fruiting condition. However, this seasonal bias likely overlooks key  
445 developmental transitions (*e.g.*, Poethig, 2013), including bud formation, bud break, leaf  
446 out, fruit development, and leaf senescence (van der Schoot *et al.*, 2014). Supporting this  
447 argument, these temporal patterns were most pronounced in NE, which experiences the  
448 harshest winter climates of the three regions. Plants collected during the winter season are  
449 almost always in dormant condition, and often lack the leaves and reproductive structures  
450 needed for taxonomic research. Collecting was also more likely during holidays and  
451 school vacations in NE and AU.

452

453 ***Trait bias***



454 In all three regions, short to medium-height species were collected more frequently than  
455 tall species (>5 m). This pattern is presumably related to the relative ease of collecting  
456 specimens from shorter, often herbaceous, species, and because reproductive materials  
457 are more accessible and potentially more abundant. Specimens of trees with woody twigs  
458 also are typically bulkier and more difficult to prepare, which may reduce their collection  
459 frequency.

460 Threatened species were also greatly under-represented in all floras. This is  
461 perhaps not surprising given their limited abundance (Palmer *et al.*, 2002) and imposed  
462 collecting restrictions (Klemens & Thorbjarnarson, 1995; Pritchard, 1996; Gibbons *et al.*,  
463 2000; Robinson, 2001). However, it is also true that collectors sometimes oversample  
464 rare or threatened species because of their higher scientific value and avoid the more  
465 common ones (Garcillán *et al.*, 2008; Garcillán & Ezcurra, 2011; Minter *et al.*, 2014).  
466 Regardless of past practices and contemporary formal restrictions, botanists now often  
467 avoid over-collection of such species by following informal guidelines and collecting  
468 plants only in areas with numerous individuals of the species (Iwanycki, 2009). Although  
469 great care in collecting rare plants is important, under-collection of rare species may lead  
470 to incorrect extinction risk assessments (*i.e.*, that the species is rarer than it actually is)  
471 and greatly limit opportunities to glean historic population and biogeographic data to  
472 guide species conservation and restoration.

473 Annuals were over-represented relative to perennials in herbarium collections in  
474 AU; the opposite was observed in SA and NE. There was also a high representation of  
475 graminoids in herbarium collections in AU and SA. This result may stem from the higher  
476 likelihood of common species being collected multiple times by different individuals or  
477 expeditions. Along these lines, much of AU is dominated by annual grasses, and the  
478 savannas of SA are populated by a variety of native and non-native perennial grasses  
479 interspersed with forbs and woody plants (Bond & Parr, 2010). New England, on the  
480 other hand, is generally forested and has an abundance of shade tolerant shrubs and  
481 perennial herbs. Graminoids are also considered harder to identify and may be avoided by  
482 non-specialists. Lianas and vines simultaneously represent the smallest proportion of  
483 growth forms and comprise the least number of specimens per species in all three floras.  
484 Such trait-based biases in botanical collections not only influence our perception of

485 species abundance and range, but can also lead to erroneous estimations of functional  
486 diversity and ecosystem services, especially for studies relying on specimen databases  
487 (Schmidt-Lebuhn *et al.*, 2013). Whether herbarium records represent true patterns of  
488 abundance and diversity remains difficult to untangle from human-mediated collecting  
489 biases. However, the large differences among the three floras in the traits of botanical  
490 collections almost certainly is reflective of genuine difference species abundance and  
491 diversity.

492

### 493 ***Phylogenetic bias***

494 Taxonomic biases in collection data have been reported previously (Hijmans *et al.*, 2000;  
495 Tobler *et al.*, 2007; Meyer *et al.*, 2016). However, our study is the first, to our knowledge,  
496 to demonstrate explicit evidence for phylogenetic bias in herbarium collections.

497 Collection efforts in all three floras were concentrated in particular clades.

498 Previous examinations of taxonomic bias (*e.g.*, Hijmans *et al.*, 2000; Tobler *et al.*,  
499 2007; Meyer *et al.*, 2016) did not use the full complement of modern phylogenetic  
500 methods that included patterns of evolutionary relatedness, and so were limited in their  
501 ability to detect details of taxonomic bias. In contrast, our phylogenetic approach not only  
502 captured taxonomic bias in favor of certain entire families (*e.g.*, Asteraceae, Cyperaceae,  
503 Poaceae, and Rosaceae in NE), but revealed that evolutionarily distinct and globally  
504 endangered species are underrepresented in herbarium records relative to more common  
505 species. Such evolutionarily distinct species, which are threatened with extinction,  
506 represent important targets for future documentation or prioritization for conservation  
507 (Isaac *et al.*, 2007). However, collecting threatened taxa requires specialized training,  
508 compliance with regulation, and awareness of actual collection needs (Minteer *et al.*,  
509 2014). Increasingly, DNA barcoding approaches, using small samples from living tissues,  
510 combined with GPS-referenced digital photography might be an avenue to document  
511 such species.

512

### 513 ***Collector bias***

514 In all three regions, a large percentage of specimens was gathered by only a few  
515 collectors (Fig. 8). Thus, the habits and preferences of a few individuals likely shaped the  
516 establishment and formation of these herbarium collections. These ‘founder effects’  
517 propagate across all the dimensions of collection bias examined above, and help us to  
518 understand past collection behavior. For example, certain collectors may focus on  
519 geographically circumscribed floristic zones, often near their place of residence,  
520 workplace, or vacation home, and sample all species found therein, whereas others may  
521 focus on collecting species of a particular clade across various regions. Professional  
522 botanists may tend to collect specimens on weekdays during any time of the year,  
523 whereas amateurs and faculty with teaching responsibilities may focus their efforts on  
524 weekends and vacation months. Those interested in function and physiology may only  
525 collect plants of certain habits or life-histories (*e.g.*, carnivorous, aquatic plants, or  
526 succulent plants). These effects would likely be compounded when associated with mega-  
527 collectors. For instance, the Harvard University Herbaria’s collection of Asian, especially  
528 woody plants, was largely built by a few collectors and dates to the early establishment of  
529 the institution, and continues to attract scholars of the flora of Asia and their collections.  
530 Investigating the historical significance and potential biases created and propagated by  
531 these early pioneers is a ripe area for future research.

532

### 533 **Future collecting**

534 To ensure that herbaria continue to be vital centers for research beyond their importance  
535 to taxonomy and systematics, herbarium directors and collectors should account for and,  
536 whenever possible, reduce biases in plant collections. Biases can be accounted for to a  
537 degree using statistical approaches (Droissart *et al.*, 2012; Feeley, 2012; Grass *et al.*,  
538 2014; Engemann *et al.*, 2015). For instance, inclusion of covariates for distances of  
539 collections from herbaria, roads, or other infrastructure (McCarthy *et al.*, 2012), using  
540 rarefaction methods to predict abundances (Schmidt-Lebuhn *et al.*, 2013), or including  
541 the collector as a variable, would improve species distribution models and associated  
542 predictions of future changes across a flora. To remedy such biases, future collecting  
543 expeditions should focus on “coldspots” of collection intensity (Hijmans *et al.*, 2000),

544 that is, places that are under-represented in collections. Although some of the coldspots  
545 we identified likely represent more inaccessible environments, they often correspond to  
546 unique ecosystems, including the Succulent Karoo of SA and the Great North Woods in  
547 northern NE that contain many species of interest. Some of these coldspots also may  
548 indicate areas where herbarium specimens have yet to be mobilized, providing additional  
549 focus for efforts to make collection data widely available. Equally important is the need  
550 to continue modern collecting in well-established “hotspots” so that there are multiple  
551 temporal benchmarks against which change can be measured. This is particularly true for  
552 non-native invasive species that have rapidly expanding distributions and vulnerable  
553 native species that have ranges that are collapsing.

554         Phylogenetic and trait biases can be alleviated by targeting collection efforts  
555 where we know species have been under-collected. Temporal bias is more difficult to  
556 address, as we cannot add to historic collections. However, we can make efforts to  
557 maintain consistent regional botanical records by conducting field surveys at regular  
558 intervals. Also, by linking multiple herbaria into larger digital databases, the temporal  
559 biases of individual herbaria can be smoothed out to some extent.

560         We acknowledge that some of the biases also may be attributed to longstanding  
561 curation practices at the herbariums themselves. As herbarium collections were amassed  
562 for qualitative floristic, taxonomic, and systematic research, duplicate specimens of  
563 common species and non-reproductive material have sometimes been discarded, sent  
564 elsewhere, or not accepted in the first place. This trend is becoming even more  
565 pronounced as many herbaria around the world are increasingly constrained by funding,  
566 labor, and space. As new uses for biological collections continue to proliferate, curation  
567 practices should also change to accommodate different avenues of research, such as  
568 climate-change biology and rare plant conservation. This will often be most effective  
569 through continued collecting of specimens to overcome past biases. And most  
570 importantly, researchers analyzing herbarium specimens in a widening array of studies  
571 needed to be aware of the biases in these collections, and apply appropriate statistical  
572 techniques.

573

574

575 **ACKNOWLEDGMENTS**

576 We thank the Harvard University Herbaria for logistic and financial support, and the  
577 virtual herbaria in the three regional floras for granting us access to their data: Australian  
578 Virtual Herbarium (<http://avh.chah.org.au>), South African National Biodiversity Institute  
579 (<http://newposa.sanbi.org/>) and the Consortium for Northeast Herbaria  
580 (<http://portal.neherbaria.org/portal/>). Digitization of most New England specimens was  
581 funded by the ADBC program of the U.S. National Science Foundation (Awards  
582 1208829, 1208835, 1208972, 1208973, 1208975, 1208989, 1209149). Special thanks to  
583 T.J. Davies, E.K. Meineke, K.M. Peterson, and K.G. Dexter for valuable discussion  
584 during the formation of this manuscript. We appreciate the constructive comments of the  
585 associate editor and three anonymous reviewers on the submitted manuscript.

586

587

588 **Author Contributions:** Conceived the project: CCD. Designed the experiment: BHD,  
589 DSP. Performed the experiments: BHD. Analyzed the data: BHD with help from  
590 DSP. Contributed reagents/materials/analysis tools: BHD, DSP, CGW, AME,  
591 CCD. Wrote the paper: BHD with significant comments and editing from all co-  
592 authors, particularly DSP, CCD, and AME.

593 **References**

594 **Abouheif E. 1999.** A method for testing the assumption of phylogenetic independence in  
595 comparative data. *Evolutionary Ecology Research* **1**: 895–909.

596 **Agostinelli C, Lund U. 2013.** *R package 'circular': Circular Statistics (version 0.4-7).*  
597 URL <https://r-forge.r-project.org/projects/circular/>

598 **APG III (Angiosperm Phylogeny Group). 2009.** An update of the angiosperm  
599 phylogeny group classification for the orders and families of flowering plants:  
600 APG III. *Botanical Journal of the Linnean Society* **161**: 105–121.

601 **AVH. 2016.** *Australia's Virtual Herbarium, Council of Heads of Australasian Herbaria,*  
602 <http://avh.chah.org.au>, accessed on 09 June 2016.

603 **Belasco WJ. 1979.** *Americans on the road, from autocamp to motel 1910-1945.*  
604 Baltimore: Johns Hopkins University Press.

605 **Bivand RS, Pebesma E, Gomez-Rubio V. 2013.** *Applied spatial data analysis with R,*  
606 *Second edition.* Springer, NY. <http://www.asdar-book.org/>

607 **Blomberg SP, Garland T, Ives AR. 2003.** Testing for phylogenetic signal in  
608 comparative data: behavioural traits are more labile. *Evolution* **57**: 717–745.

609 **Bond WJ, Parr CL. 2010.** Beyond the forest edge: ecology, diversity and conservation  
610 of the grassy biomes. *Biological Conservation* **143**: 2395–2404.

611 **Botkin BA. 1968.** Automobile humor: from the horseless carriage to the compact car.  
612 *The Journal of Popular Culture* **I**: 395–402.

613 **Boyle B, Hopkins N, Lu Z, Garay JAR, Mozzherin D, Rees T, Matasci N, Narro ML,**  
614 **Piel WH, Mckay SJ, et al. 2013.** The taxonomic name resolution service: an  
615 online tool for automated standardization of plant names. *BMC Bioinformatics* **14**:  
616 16.

617 **Ceballos G, Ehrlich PR. 2006.** Global mammal distributions, biodiversity hotspots, and  
618 conservation. *Proceedings of the National Academy of Sciences USA* **103**: 19374–  
619 19379.

620 **CIESIN. 2016.** *Center for International Earth Science Information Network, Columbia*  
621 *University.* Gridded Population of the World, Version 4 (GPWv4): Population  
622 Density. Palisades, NY: NASA Socioeconomic Data and Applications Center  
623 (SEDAC). <http://dx.doi.org/10.7927/H4NP22DQ>. Accessed 29 August 2016.

- 624 **CNH. 2016.** Consortium of Northeastern Herbaria. <http://portal.neherbaria.org/portal/>
- 625 **Cotterill FPD, Hustler CW, Broadley DG. 1994.** Systematics and biodiversity. *Trends*  
626 *in Ecology & Evolution* **9**: 228.
- 627 **Cotterill FPD. 1995.** Systematics, biological knowledge and environmental conservation.  
628 *Biodiversity and Conservation* **4**: 183–205.
- 629 **Dalton R. 2003.** Natural history collections in crisis as funding is slashed. *Nature* **423**:  
630 6940.
- 631 **Daru BH, Elliott TL, Park DS, Davies TJ. 2017.** Understanding the processes  
632 underpinning patterns of phylogenetic regionalization. *Trends in Ecology &*  
633 *Evolution* doi: 10.1016/j.tree.2017.08.013
- 634 **Daru BH, Van der Bank M, Davies TJ. 2015.** Spatial incongruence among hotspots  
635 and complementary areas of tree diversity in southern Africa. *Diversity and*  
636 *Distributions* **21**: 769–780.
- 637 **Davies TJ, Kraft NJB, Salamin N, Wolkovich EM. 2012.** Incompletely resolved  
638 phylogenetic trees inflate estimates of phylogenetic conservatism. *Ecology* **93**:  
639 242–247.
- 640 **Davis CC, Willis CG, Connolly B, Kelly C, Ellison AM. 2015.** Herbarium records are  
641 reliable sources of phenological change driven by climate and provide novel  
642 insights into species' phenological cueing mechanisms. *American Journal of*  
643 *Botany* **102**: 1599–1609.
- 644 **Droissart V, Hardy OJ, Sonké B, Dahdouh-Guebas F, Stévant T. 2012.** Subsampling  
645 herbarium collections to assess geographic diversity gradients: A case study with  
646 endemic Orchidaceae and Rubiaceae in Cameroon. *Biotropica* **44**: 44–52.
- 647 **Edwards EJ, de Vos JM, Donoghue MJ. 2015.** Doubtful pathways to cold tolerance in  
648 plants. *Nature* **521**: E5–E6.
- 649 **Edwards JL, Lane MA, Nielsen ES. 2000.** Interoperability of biodiversity databases:  
650 biodiversity information on every desktop. *Science* **289**: 2312–2314.
- 651 **Engemann K, Enquist BJ, Sandel B, Boyle B, Jørgensen PM, Morueta-Holme N,**  
652 **Peet RK, Violle C, Svenning J-C. 2015.** Limited sampling hampers “big data”  
653 estimation of species richness in a tropical biodiversity hotspot. *Ecology and*  
654 *Evolution* **5**: 807–820.

- 655 **ESRI 1992.** *Environmental Systems Research Institute, Digital chart of the world, 1:1M.*  
656 Environmental Systems Research Institute, Inc., Redlands, California.
- 657 **Everill PH, Primack RB, Ellwood EE, Melaas EK. 2014.** Determining past leaf-out  
658 times of New England's deciduous forests from herbarium specimens. *American*  
659 *Journal of Botany* **101**: 1–8.
- 660 **Feeley KJ. 2012.** Distributional migrations, expansions, and contractions of tropical plant  
661 species as revealed in dated herbarium records. *Global Change Biology* **18**: 1335–  
662 1341.
- 663 **Ferreira SLA, Harmse AC. 2000.** Crime and tourism in South Africa: international  
664 tourists perception and risk. *South African Geographical Journal* **82**: 80–85.
- 665 **Forman RTT, Alexander LE. 1998.** Roads and their major ecological effects. *Annual*  
666 *Review of Ecology and Systematics* **29**: 207–31
- 667 **Fortune S. 1992.** Voronoi diagrams and Delaunay triangulations. *Computing in*  
668 *Euclidean Geometry* **1**: 193–233.
- 669 **Funk V. 2003.** The importance of herbaria. *Plant Science Bulletin* **49**: 94–95.
- 670 **Funk VA, Morin N. 2000.** A survey of the herbaria of the southeast United States. *Sida,*  
671 *Botanical Miscellany* **18**: 35–52.
- 672 **Funk VA, Richardson K. 2002.** Biological specimen data in biodiversity studies: use it  
673 or lose it. *Systematic Biology* **51**: 303–316.
- 674 **GADM. 2015.** *Global Administrative Areas*, version 2.8 ([www.gadm.org](http://www.gadm.org)).
- 675 **Garcillán PP, Ezcurra E, Vega E. 2008.** Guadalupe Island: Lost paradise recovered?  
676 Overgrazing impact on extinction in a remote oceanic island as estimated through  
677 accumulation functions. *Biodiversity and Conservation* **17**: 1613–1625.
- 678 **Garcillán PP, Ezcurra E. 2011.** Sampling procedures and species estimation: Testing  
679 the effectiveness of herbarium data against vegetation sampling in an oceanic  
680 island. *Journal of Vegetation Science* **22**: 273–280.
- 681 **Geri F, Lastrucci L, Viciani D, Foggi B, Ferretti G, Maccherini S, Bonini I, Amici V,**  
682 **Chiarucci A. 2013.** Mapping patterns of ferns species richness through the use of  
683 herbarium data. *Biodiversity and Conservation* **22**: 1679–1690.



- 684 **Gibbons JW, Scott DE, Ryan T, Buhlmann K, Tuberville T, Greene J, Mills T,**  
685 **Leiden Y, Poppy S, Winne C *et al.* 2000.** The global decline of reptiles, déj·vu  
686 amphibians. *BioScience* **50**: 653–666.
- 687 **Grass A, Tremetsberger K, Hössinger R, Bernhardt K. 2014.** Change of species and  
688 habitat diversity in the Pannonian region of eastern Lower Austria over 170 years:  
689 Using herbarium records as a witness. *Natural Resources* **5**: 583–596.
- 690 **Greve M, Lykke AM, Fagg CW, Gereau RE, Lewis GP, Marchant R, Marshall AR,**  
691 **Ndayishimiye J, Bogaert J, Svenning JC. 2016.** Realising the potential of  
692 herbarium records for conservation biology. *South African Journal of Botany* **105**:  
693 317–323.
- 694 **Griffith EH, Sauer JR, Royle JA. 2010.** Traffic effects on bird counts on North  
695 American breeding bird survey routes. *Auk* **127**: 387–393.
- 696 **Gutzwiller KJ, Flather CH. 2011.** Wetland features and landscape context predict the  
697 risk of wetland habitat loss. *Ecological Applications* **21**: 968–982
- 698 **Hart R, Salick J, Ranjitkar S, Xu J. 2014.** Herbarium specimens show contrasting  
699 phenological responses to Himalayan climate. *Proceedings of the National*  
700 *Academy of Sciences USA* **111**: 10615–10619.
- 701 **Hijmans RJ, Garrett KA, Huaman Z, Zhang DP, Schreuder M, Bonierbale M. 2000.**  
702 Assessing the geographic representation of genebank collections: the case of the  
703 Bolivian wild potatoes. *Conservation Biology* **14**: 1755–1765.
- 704 **Hijmans RJ. 2015.** *geosphere: Spherical Trigonometry*. R package version 1.4-3.  
705 <http://CRAN.R-project.org/package=geosphere>
- 706 **Hortal J, Lobo JM, Jiménez-Valverde A. 2007.** Limitations of biodiversity databases:  
707 case study on seed-plant diversity in tenerife, canary islands. *Conservation*  
708 *Biology* **21**: 853–863.
- 709 **Hui C, Shuang-cheng L, Yi-li Z. 2003.** Impact of road construction on vegetation  
710 alongside Qinghai-Xizang highway and railway. *Chinese Geographical Science*  
711 **13**: 340–346.
- 712 **Isaac NJ, Turvey ST, Collen B, Waterman C, Baillie JE. 2007.** Mammals on the  
713 EDGE: conservation priorities based on threat and phylogeny. *PLoS ONE* **2**: e296.

- 714 **Iwanycki N. 2009.** *Guidelines for collecting herbarium specimens of vascular plants.*  
715 Royal Botanical Gardens Canada, Hamilton, Canada.
- 716 **Jombart T, Dray S. 2008.** adephylo: exploratory analyses for the phylogenetic  
717 comparative method. *Bioinformatics* **26**: 1907–1909.
- 718 **Klemens MW, Thorbjarnarson JB. 1995.** Reptiles as a food resource. *Biodiversity and*  
719 *Conservation* **4**: 281–298.
- 720 **Lavoie C. 2013.** Biological collections in an ever changing world: Herbaria as tools for  
721 biogeographical and environmental studies. *Perspectives in Plant Ecology,*  
722 *Evolution and Systematics* **15**: 68–76.
- 723 **le Roux MM, Wilkin P, Balkwill K, Boatwright JS, Bytebier B, Filer D, Klak C,**  
724 **Klopper RR, Koekemoer M, Livermore L et al. 2017.** Producing a plant  
725 diversity portal for South Africa. *Taxon* **66**: 421–431.
- 726 **Lees DC, Lack HW, Rougerie R, Hernandez-Lopez A, Raus T, Avtzis ND, Augustin**  
727 **S, Lopez-Vaamonde C. 2011.** Tracking origins of invasive herbivores through  
728 herbaria and archival DNA: the case of the horse-chestnut leaf miner. *Frontiers in*  
729 *Ecology and the Environment* **9**: 322–328.
- 730 **Lemanski C. 2004.** A new apartheid? The spatial implications of fear of crime in Cape  
731 Town, South Africa. *Environment & Urbanization* **16**: 101–111.
- 732 **Leuner B. 2007.** *Migration, multiculturalism and language maintenance in Australia.*  
733 Peter Lang, Oxford.
- 734 **Li Y, Yu J, Ning K, Du S, Han G, Qu F, Wang G, Fu Y, Zhan C. 2014.** Ecological  
735 effects of roads on the plant diversity of coastal wetland in the Yellow River Delta.  
736 *The Scientific World Journal* **2014**: 952051.
- 737 **McCarthy KP, Fletcher JR RJ, Rota CT, Hutto RL. 2012.** Predicting species  
738 distributions from samples collected along roadsides. *Conservation Biology* **26**:  
739 68–77.
- 740 **Meyer C, Weigelt P, Kreft H. 2016.** Multidimensional biases, gaps and uncertainties in  
741 global plant occurrence information. *Ecology Letters* **19**: 992–1006.
- 742 **Miller-Rushing A, Primack R, Mukunda S. 2006.** Photographs and herbarium  
743 specimens as tools to document phenological changes in response to global  
744 warming. *American Journal of Botany* **93**: 1667–1674.

- 745 **Minteer BA, Collins JP, Love KE, Puschendorf R. 2014.** Avoiding (Re)Extinction".  
746 *Science* **344**: 260-261.
- 747 **Moerman DE, Estabrook GF. 2006.** The botanist effect: counties with maximal species  
748 richness tend to be home to universities and botanists. *Journal of Biogeography*  
749 **33**: 1969–1974.
- 750 **Musgrave T, Gardner C, Musgrave W. 1999.** *The plant hunters. Two hundred years of*  
751 *adventure and discovery*. Seven Dials, London, UK.
- 752 **NEBC. 1899.** Editorial announcement. *Rhodora* **1**: 1–2
- 753 **Newbold T. 2010.** Applications and limitations of museum data for conservation and  
754 ecology, with particular attention to species distribution models. *Progress in*  
755 *Physical Geography* **34**: 3–22.
- 756 **Norris WR, Lewis DQ, Widrlechner MP, Thompson JD, Pope RO. 2001.** Lessons  
757 from an inventory of the Ames, Iowa, flora (1859–2000). *Journal of the Iowa*  
758 *Academy of Science* **108**: 34–63.
- 759 **Nualart N, Ibáñez N, Soriano I, López-Pujol J. 2017.** Assessing the relevance of  
760 herbarium collections as tools for conservation biology. *Botanical Review* **83**:  
761 303–325.
- 762 **Orme CD, Davies RG, Burgess M, Eigenbrod F, Pickup N, Olson VA, Webster AJ,**  
763 **Ding TS, Rasmussen PC, Ridgely RS, et al. 2005.** Global hotspots of species  
764 richness are not congruent with endemism or threat. *Nature* **436**: 1016–1019.
- 765 **Orme D, Freckleton R, Thomas G, Petzoldt T, Fritz S, Isaac N, Pearse W. 2012.**  
766 *caper: Comparative Analyses of Phylogenetics and Evolution in R*. R package  
767 version 0.5. [http://CRAN.R-project.org/package = caper](http://CRAN.R-project.org/package=caper).
- 768 **Pagel M. 1999.** Inferring the historical patterns of biological evolution. *Nature* **401**: 877–  
769 884.
- 770 **Palmer MW, Earls PG, Hoagland BW, White PS, Wohlgemuth T 2002.** Quantitative  
771 tools for perfecting species list. *Environmetrics* **13**: 121–137.
- 772 **Pautasso M, McKinney ML. 2007.** The Botanist Effect Revisited: Plant Species  
773 Richness, County Area, and Human Population Size in the United States.  
774 *Conservation Biology* **21**: 1333–1340.

- 775 **Poethig, RS. 2013.** Vegetative phase change and shoot maturation in plants. *Current*  
776 *Topics in Developmental Biology* **105**: 125–152.
- 777 **Prather LA, Alvarez-Fuentes O, Mayfield MH, Ferguson CJ. 2004a.** The decline of  
778 plant collecting in the United States: a threat to the infrastructure of biodiversity  
779 studies. *Systematic Botany* **29**: 15–28.
- 780 **Prather LA, Alvarez-Fuentes O, Mayfield MH, Ferguson CJ. 2004b.** Implications of  
781 the decline in plant collecting for systematic and floristic research. *Systematic*  
782 *Botany* **29**: 216–220.
- 783 **Price CA. 1998.** Post-war immigration: 1945-1998. *Journal of the Australian Population*  
784 *Association* **15**: 17.
- 785 **Pritchard PCH. 1996.** *The Galápagos tortoises: nomenclatural and survival status.*  
786 Chelonian Research Foundation in association with Conservation International  
787 and Chelonia Institute, Lunenburg, MA.
- 788 **Pyke GH, Ehrlich PR. 2010.** Biological collections and ecological/environmental  
789 research: a review, some observations and a look to the future. *Biological Reviews*  
790 **5**: 247–266.
- 791 **Redding DW, Mooers AØ. 2006.** Incorporating evolutionary measures into conservation  
792 prioritization. *Conservation Biology* **20**: 1670–1678.
- 793 **Revell LJ. 2012.** phytools: An R package for phylogenetic comparative biology (and  
794 other things). *Methods in Ecology and Evolution* **3**: 217–223.
- 795 **Rich TCG, Woodruff ER. 1992.** Recording bias in botanical surveys. *Watsonia* **19**: 73–  
796 95.
- 797 **Rivers-Moore NA, Cowden C. 2012.** Regional prediction of wetland degradation in  
798 South Africa. *Wetlands Ecology and Management* **20**: 491–502.
- 799 **Robinson JG. 2001.** Using ‘sustainable use’ approaches to conserve exploited  
800 populations. In: Reynolds JD, Mace GM, Redford KH, Robinson JG, eds.  
801 *Conservation of exploited species.* Cambridge: Cambridge University Press, 485–  
802 498.
- 803 **Rudin SM, Murray DW, Whitfeld TJS. 2017.** Retrospective analysis of heavy metal  
804 contamination in Rhode Island based on old and new herbarium specimens.  
805 *Applications in Plant Sciences* **5**: 1–13.

- 806 **SANBI. 2016.** *South African National Biodiversity Institute*. Botanical Database of  
807 Southern Africa (BODATSA), <http://newposa.sanbi.org/>, accessed on 22 July  
808 2016.
- 809 **Schaefer H, Hardy OJ, Silva L, Barraclough TG, Savolainen V. 2011.** Testing  
810 Darwin's naturalization hypothesis in the Azores. *Ecology Letters* **14**: 389–396.
- 811 **Schmidt-Lebuhn AN, Knerr NJ, Kessler M. 2013.** Non-geographic collecting biases in  
812 herbarium specimens of Australian daisies (Asteraceae). *Biodiversity and*  
813 *Conservation* **22**: 905–919.
- 814 **Schorn C, Weber E, Bernardos R, Hopkins C, Davis CC. 2016.** The New England  
815 Vascular Plants Project: 295,000 specimens and counting. *Rhodora* **118**: 324–325.
- 816 **Staats M, Erkens RHJ, van de Vossen B, Wieringa JJ, Kraaijeveld K, Stielow B,**  
817 **Geml J, Richardson JE, Bakker FT. 2013.** Genomic treasure troves: complete  
818 genome sequencing of herbarium and insect museum specimens. *PLoS ONE* **8**:  
819 e69189.
- 820 **Stropp J, Ladle RJ, Malhado ACM, Hortal J, Gaffuri J, Temperley, WH, Skøien JO.**  
821 **Mayaux, P. 2016.** Mapping ignorance: 300 years of collecting flowering plants in  
822 Africa. *Global Ecology and Biogeography* **25**: 1085–1096.
- 823 **Syfert MM, Smith MJ, Coomes DA. 2013.** The effects of sampling bias and model  
824 complexity on the predictive performance of MaxEnt species distribution models.  
825 *PLoS ONE* **8**: e55158.
- 826 **Thiers B. 2016.** *Index Herbariorum: A global directory of public herbaria and*  
827 *associated staff*. New York Botanical Garden's Virtual Herbarium.  
828 <http://sweetgum.nybg.org/science/ih/>, accessed on 29 September 2016.
- 829 **Tobler M, Honorio E, Janovec J, Reynel C. 2007.** Implications of collection patterns of  
830 botanical specimens on their usefulness for conservation planning: an example of  
831 two neotropical plant families (Moraceae and Myristicaceae) in Peru. *Biodiversity*  
832 *and Conservation* **16**: 659–677
- 833 **van der Schoot C, Paul LK, Rinne PLH. 2014.** The embryonic shoot: a lifeline through  
834 winter. *Journal of Experimental Botany* **65**: 1699–1712.

- 835 **Webb CO, Ackerly DD, Kembel SW. 2008.** PHYLOCOM: software for the analysis of  
836 phylogenetic community structure and trait evolution. *Bioinformatics* **24**: 2098–  
837 2100.
- 838 **Webb CO, Ackerly DD, McPeck MA, Donoghue MJ. 2002.** Phylogenies and  
839 community ecology. *Annual Review of Ecology and Systematics* **33**: 475–505.
- 840 **Webb CO, Donoghue MJ. 2005.** Phylomatic: tree assembly for applied phylogenetics.  
841 *Molecular Ecology Notes* **5**: 181–183.
- 842 **Whittle T. 1970.** *The Plant Hunters*. Heinemann, London.
- 843 **Willis CG, Ellwood ER, Primack RB, Davis CC, Pearson KD, Gallinato AS, Yost**  
844 **JM, Nelson G, Mazer SJ, Rossington NL et al. 2017a.** Old plants, new tricks:  
845 phenological research using herbarium specimens. *Trends in Ecology & Evolution*  
846 **32**: 531–546.
- 847 **Willis CG, Law E, Williams AC, Franzone BF, Bernardos R, Brun L, Hopkins C,**  
848 **Schorn C, Weber E, Parks DS et al. 2017b.** CrowdCurio: an online  
849 crowdsourcing platform to facilitate climate change studies using herbarium  
850 specimens. *New Phytologist* **215**: 479–488.
- 851 **Wolf A, Anderegg WRL, Ryan SJ, Christensen J. 2011.** Robust detection of plant  
852 species distribution shifts under biased sampling regimes. *Ecosphere* **2**: 115.
- 853 **Wolkovich EM, Davies TJ, Schaefer H, Cleland EE, Cook BI, Travers SE, Willis**  
854 **CG, Davis CC. 2013.** Temperature-dependent shifts in phenology contribute to  
855 the success of exotic species with climate change. *American Journal of Botany*  
856 **100**: 1407–1421.
- 857 **Yessoufou K, Daru BH, Davies TJ. 2012.** Phylogenetic patterns of extinction risk in the  
858 Eastern Arc ecosystems, an African biodiversity hotspot. *PLoS ONE* **7**: e47082.
- 859 **Zanne AE, Tank DC, Cornwell WK, Eastman JM, Smith SA, FitzJohn RG,**  
860 **McGlenn DJ, O'Meara BC, Moles AT, Reich PB, et al. 2014.** Three keys to the  
861 radiation of angiosperms into freezing environments. *Nature* **506**: 89–92.
- 862
- 863

864 **Supporting Information**

865 Additional supporting information may be found in the online version of this article.

866

867 **Fig. S1** Analytical workflow representing the different steps in the development of this  
868 study from data compilation, collation, to statistical analysis.

869

870 **Fig. S2** Relationships of the number of specimens collected per species with number of  
871 species collected in each flora for Australia (left), South Africa (middle), and New  
872 England of the USA (right).

873

874 **Table S1** Results of the tests of phylogenetic signal in the number of specimens collected  
875 per species.

876

877

878 **LEGENDS TO FIGURES**

879 **Fig. 1** Spatial bias in herbarium collections. Geographic distribution of herbarium  
880 collecting activity depicting the spatial variation in sampling effort using Delaunay  
881 polygon tiles for (a) Australia (857,245 locales), (b) South Africa (n = 61,130 locales),  
882 and (c) New England (n = 130, 374 locales). Hotspots (red) and coldspots (blue) of  
883 herbarium sampling within quarter degree grids for (d) Australia, (e) South Africa and (f)  
884 New England. The hotspots and coldspots are the top and lowest 2.5% quantiles  
885 respectively of the number of specimens per locale.

886

887 **Fig. 2** Comparison of geographic sampling bias of herbarium records in relation to (a) the  
888 minimum distance to roads; (b) minimum distance to herbaria; and (c) regional altitudes  
889 at sampling locales. Black lines in (a) and (b) correspond to sampling locales and red  
890 indicates an equal number of random points generated 1000 times. Dark grey shading in  
891 (c) corresponds to sampling locales in relation to the regional altitudes, *i.e.*, all other  
892 altitudes (in red) for all three floras, Australia (left), South Africa (middle) and New  
893 England (right). Dotted line in (c) indicates altitude at 500 m above sea level.

894

895 **Fig. 3** Timeline of herbarium specimen collection density in relation to major historical  
896 events in time (indicated in red text) for the three floras: Australia, South Africa and New  
897 England. Analysis of phylogenetic structure through time by binning sequences of  
898 collection dates into decades and testing for overdispersion *vs.* clustering, are indicated in  
899 black font. The red trend line indicates the gross domestic product of each region. NRI,  
900 net relatedness index; NTI, nearest taxon index.

901

902 **Fig. 4** Temporal biases in herbarium collections. (a) Comparison of density plots of  
903 collection dates by seasons of the year of herbarium records (blue line) with the dates  
904 spanning the entire duration of collection (red line); blue lines outside the red lines  
905 indicate over-collecting at a particular time of year, and (b) Distribution of collection  
906 dates by days of the week for the three floras. Australia (n = 4,579,321 collection dates),  
907 South Africa (n = 771,991 collection dates), and New England (n = 562,587 collection  
908 dates).



909

910 **Fig. 5** Assessment of bias in plant traits: (a) growth duration; (b) growth form; (c) height;  
911 and (d) extinction risk for the floras of Australia (left pane), South Africa (middle pane)  
912 and New England (right pane). Error bars in (a), (b), and (d) represent +/- SE.

913

914 **Fig. 6** Distribution of phylogenetic bias, the tendency of closely related species to be  
915 similarly collected in herbarium records for three floras: (a) Australia; (b) South Africa;  
916 and (c) New England. Collecting effort is not phylogenetically random, but tends to be  
917 clustered in few selected lineages. The color scales correspond to  $\log_{10}$  numbers of  
918 specimens per species and ranges from red (low number of specimens per species) to blue  
919 (high number of specimens per species).

920

921 **Fig. 7** Phylogenetic bias in collection frequency for exemplar families in New England  
922 flora. Phylogenetic bias is indicated by significant phylogenetic signal in at least one of  
923 three metrics (Abouheif's  $C_{\text{mean}}$ , Blomberg's  $K$  and Pagel's  $\lambda$ ). The color bar illustrates  
924 values within families:  $\log_{10}$  numbers of specimens per species and ranges from red (low  
925 number of specimens per species) to blue (high number of specimens per species). \*\* $P <$   
926  $0.001$ ; \* $P < 0.01$ ; NS  $P > 0.05$

927

928 **Fig. 8** Collector bias in herbarium collections. The number of herbarium specimens  
929 amassed per collector for three regional floras in (a) Australia; (b) South Africa; and (c)  
930 New England. The top five collectors in each flora are highlighted in red. Numbers  
931 within parentheses correspond to lifespans of the collectors, with collectors that have died  
932 highlighted in red and currently living ones in black.

933

934

935 **Table 1.** Model coefficients for multiple regressions of collecting effort in the number of specimens collected per locality.

AUSTRALIA	Predictors (log <sub>10</sub> -transformed)	Percentage of variance explained (%)	P values	Model adjusted R <sup>2</sup>	Model slope	Model intercept
	Distance to roads	0.14	0.001	0.4571	-0.02	11.45
	Distance to herbaria	45.03	0.001		-0.89	
	Human population density	0.50	0.001		0.11	
	Altitude	0.041	0.001		-0.046	
SOUTH AFRICA	Predictors (log <sub>10</sub> -transformed)	Percentage of variance explained (%)	P values	Model adjusted R <sup>2</sup>	Model slope	Model intercept
	Distance to roads	0.00001	0.0003	0.3075	-0.011	11.33
	Distance to herbaria	29.13	0.001		-0.73	
	Human population density	0.0009	0.001		-0.03	

	Altitude	1.62	0.001		-0.15	
NEW ENGLAND	Predictors (log <sub>10</sub> -transformed)	Percentage of variance explained (%)	P values	Model adjusted R <sup>2</sup>	Model slope	Model intercept
	Distance to roads	0.07	0.0009	0.17	0.13	7.03
	Distance to herbaria	12.3	0.001		-0.87	
	Human population density	4.68	0.001		0.30	
	Altitude	0.04	0.001		0.046	

936

937

938

939 **Table 2.** Results of the tests of phylogenetic signal in the number of specimens collected per species using three methods (Abouheif's  
 940  $C_{\text{mean}}$ , Blomberg's K and Pagel's  $\lambda$ ). Phylogenetic data is derived from Zanne *et al.* (2014). All tests are based on 1000 randomizations.  
 941 **\*\*P < 0.001; <sup>NS</sup>P > 0.05**

	Australia (n = 5814 species)	South Africa (n = 3568 species)	New England (n = 4269 species)
Abouheif's $C_{\text{mean}}$	0.12**	0.15**	0.12**
Blomberg's K	0.00085 <sup>NS</sup>	0.0013 <sup>NS</sup>	0.0030 <sup>NS</sup>
Pagel's $\lambda$	0.18**	0.32**	0.29**

942

943

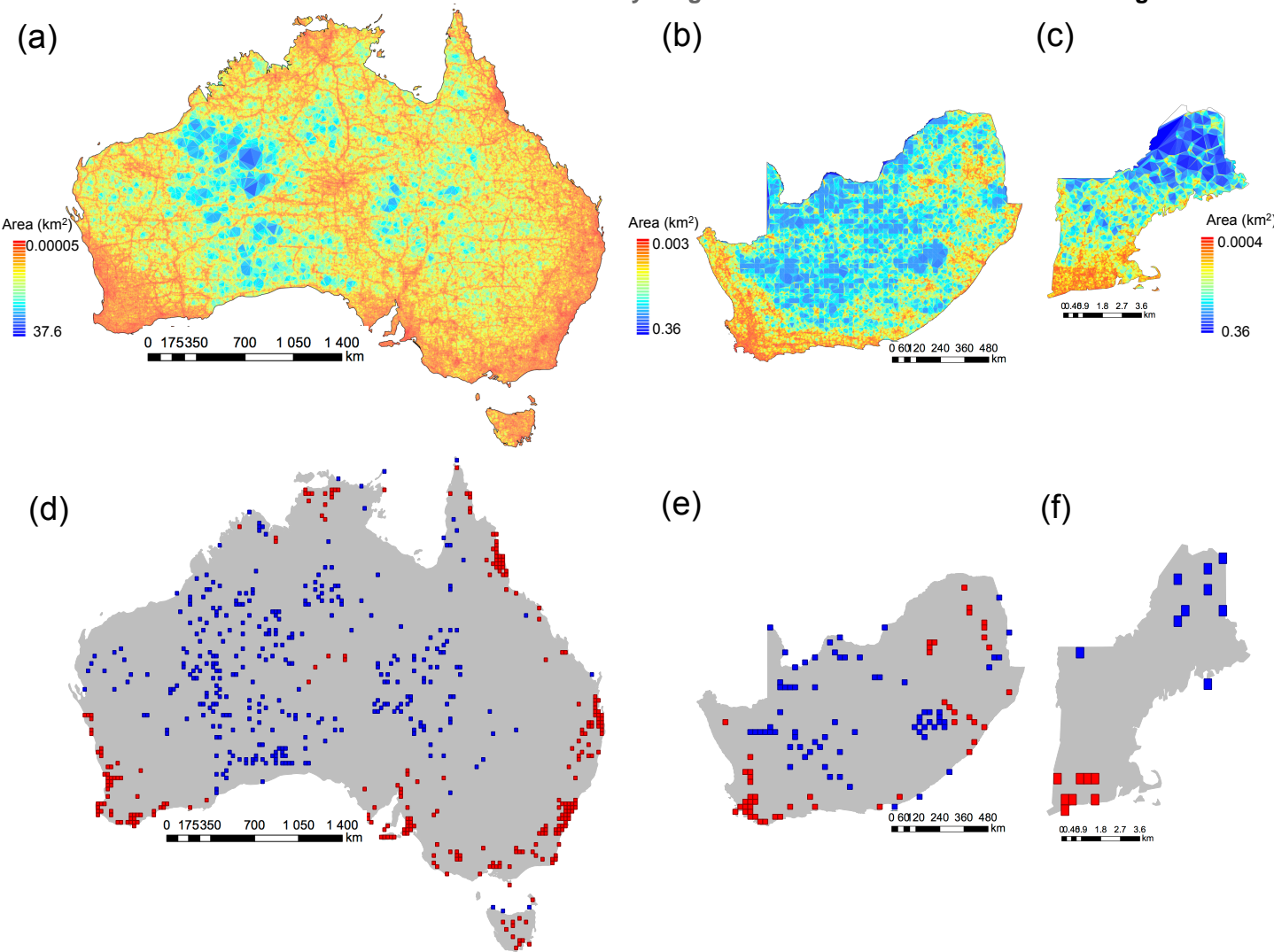
944

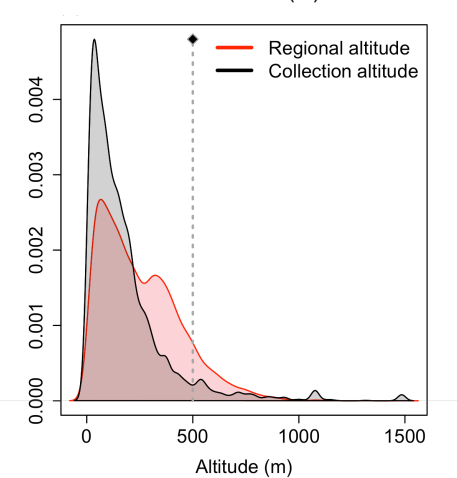
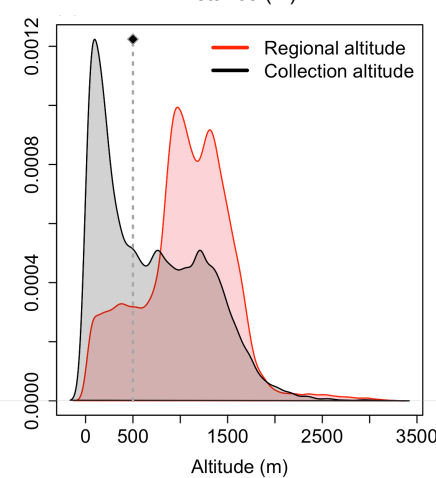
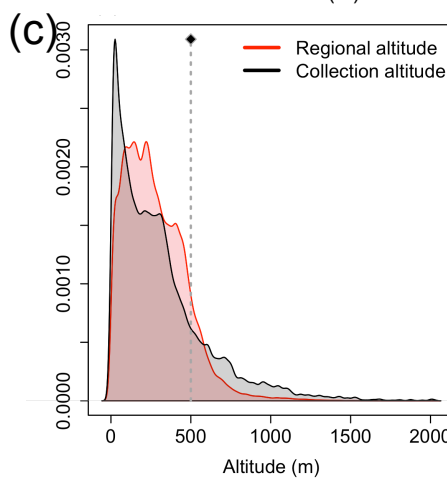
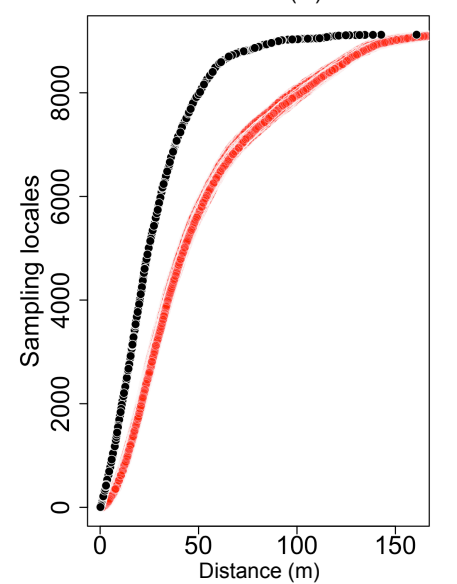
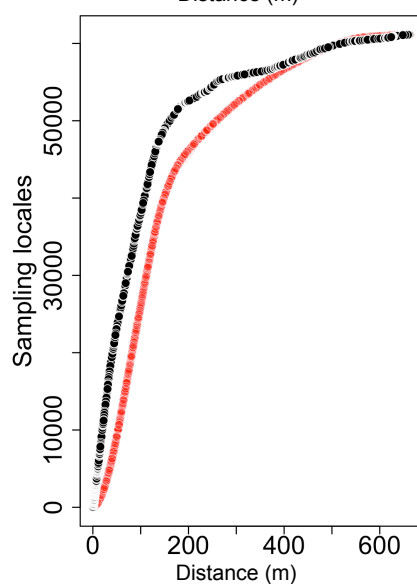
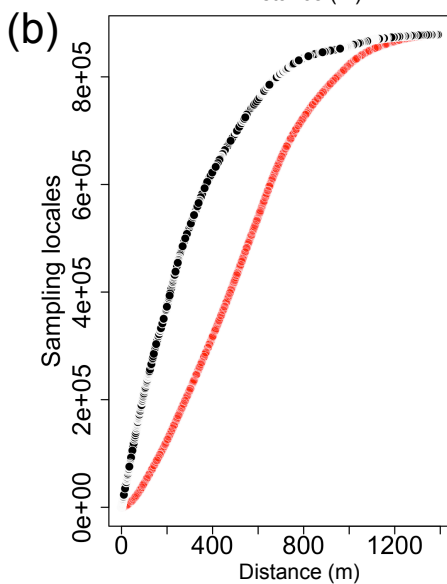
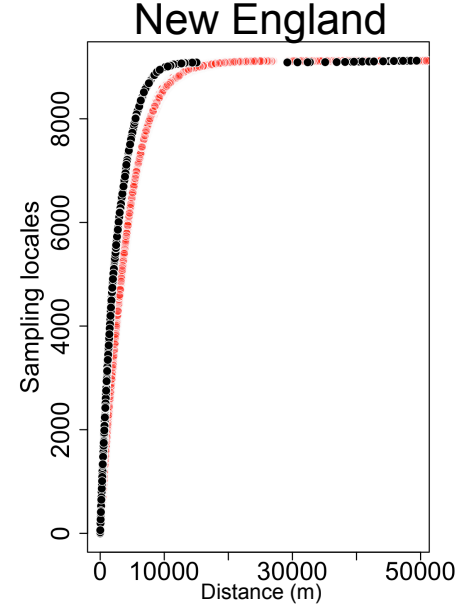
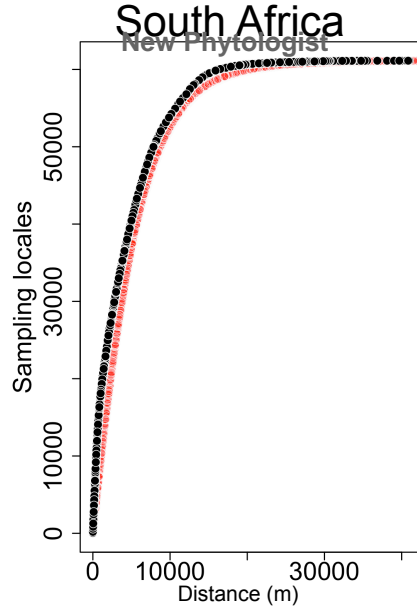
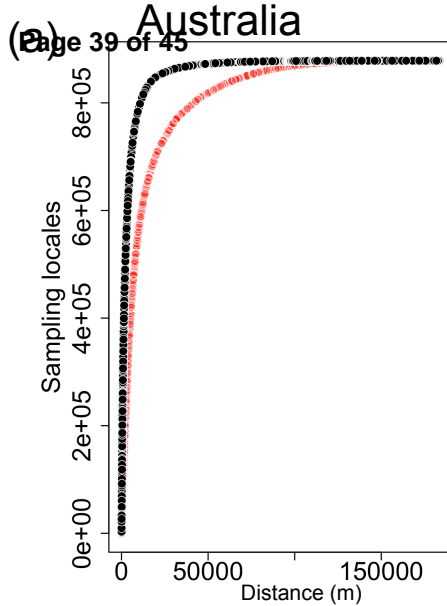
945 **Table 3** Multiple regressions of phylogenetic generalized least squares of collecting effort (frequency) of herbarium specimens with  
 946 phylogenetic metrics of species uniqueness. BL, terminal branch length; ED, evolutionary distinctiveness; EDGE, evolutionary  
 947 distinctiveness and global endangerment.

Australia	Predictors (log <sub>10</sub> -transformed)	Percentage of variance explained (%)	P values	Model adjusted R <sup>2</sup>	Model slope	Model intercept
	BL	1.36	0.7	0.049	0.035	4.37
	ED	0.2	0.008		0.44	
	EDGE	3.75	<0.001		-1.23	
South Africa	Predictors (log <sub>10</sub> -transformed)	Percentage of variance explained (%)	P values	Model adjusted R <sup>2</sup>	Model slope	Model intercept
	BL	0.47	0.3	0.09	-0.063	3.63
	ED	0.000015	0.001		0.63	
	EDGE	8.89	<0.001		-1.3	
New England	Predictors (log <sub>10</sub> -transformed)	Percentage of variance explained (%)	P values	Model adjusted R <sup>2</sup>	Model slope	Model intercept

		explained (%)				
	BL	0.09	0.94	1.70E-02	-0.0052	3.89
	ED	0.054	0.0045		0.79	
	EDGE	1.87	<0.001		-2.28	

948



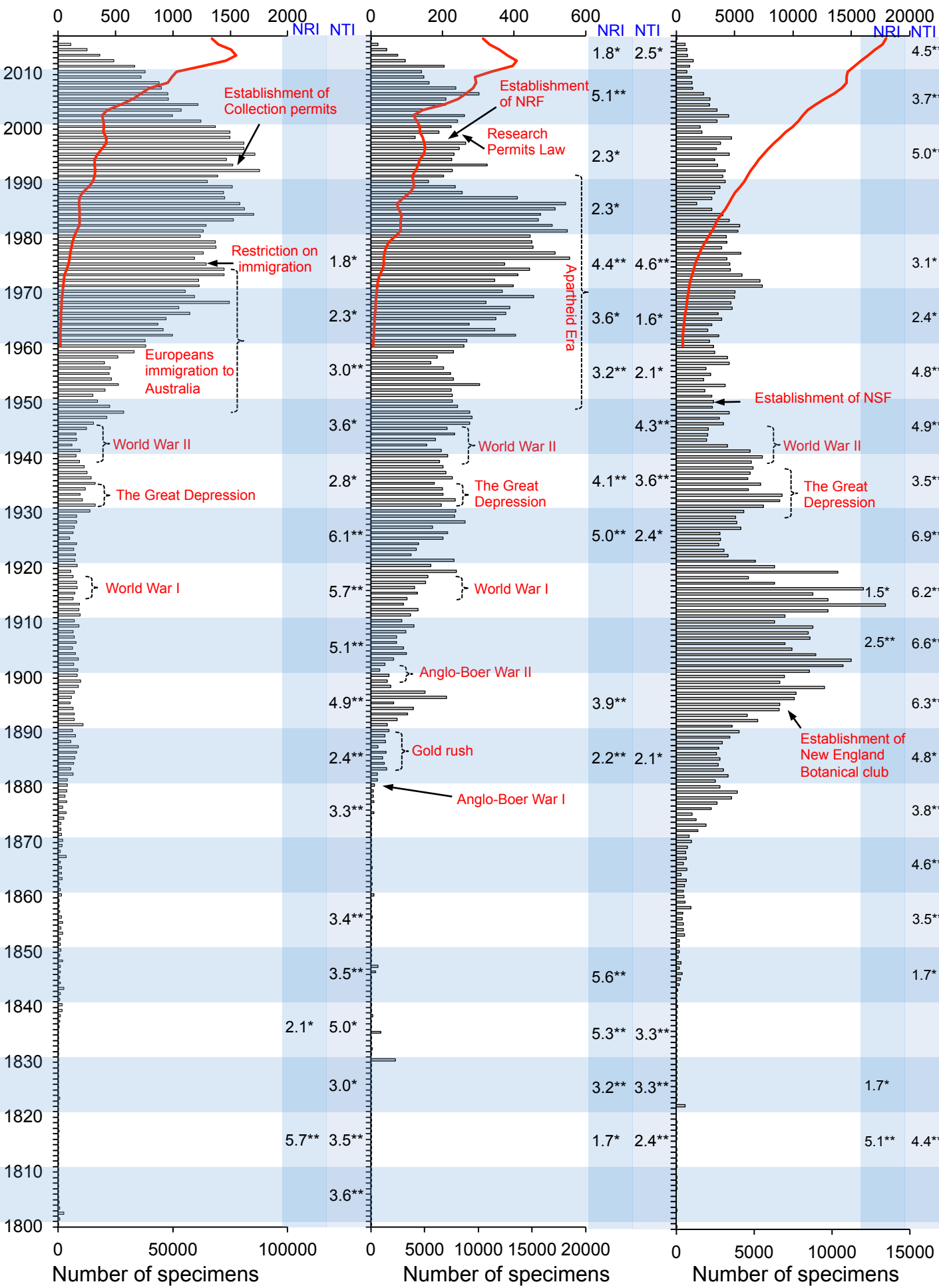




GDP (× US\$10<sup>9</sup>)

GDP (× US\$10<sup>9</sup>)

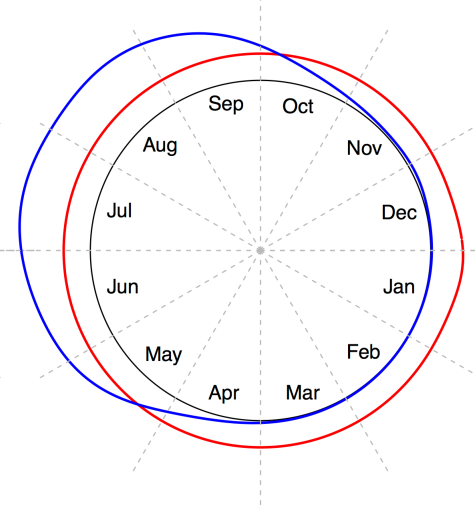
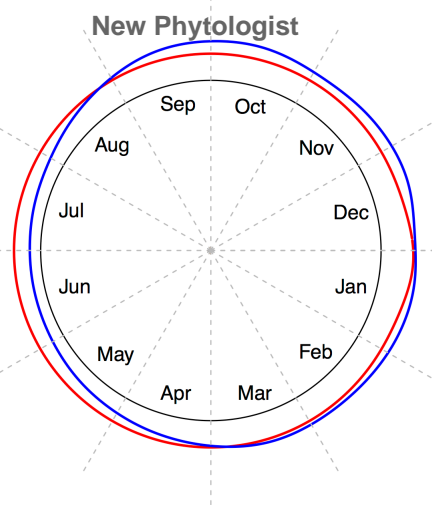
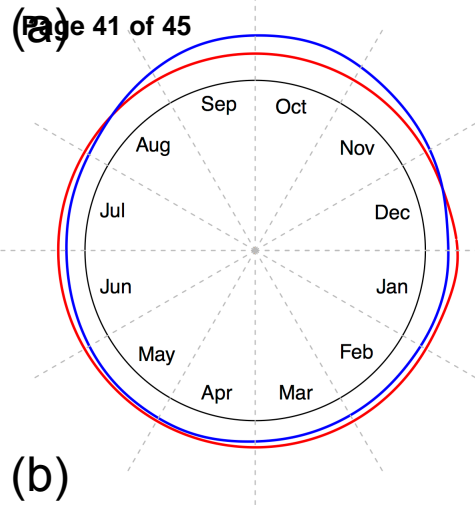
GDP (× US\$10<sup>9</sup>)



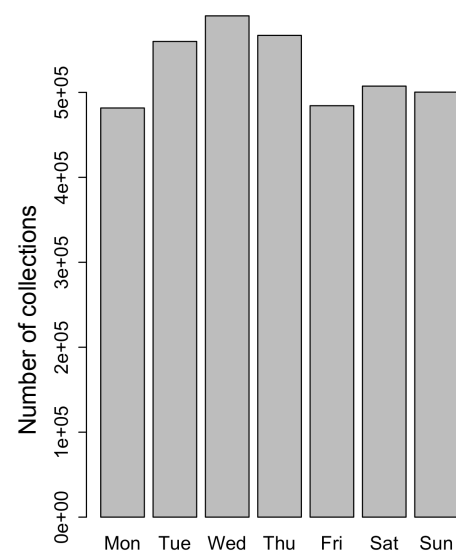
Australia

South Africa

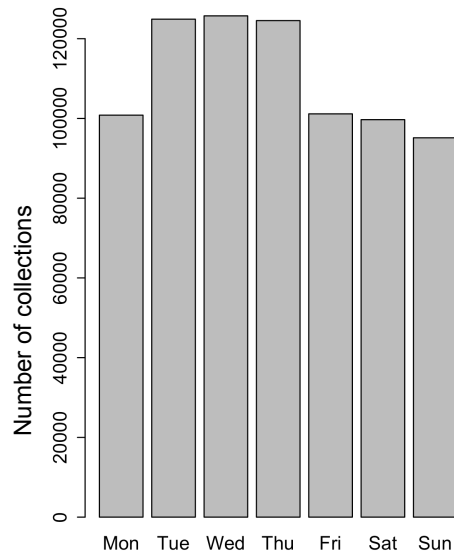
New England



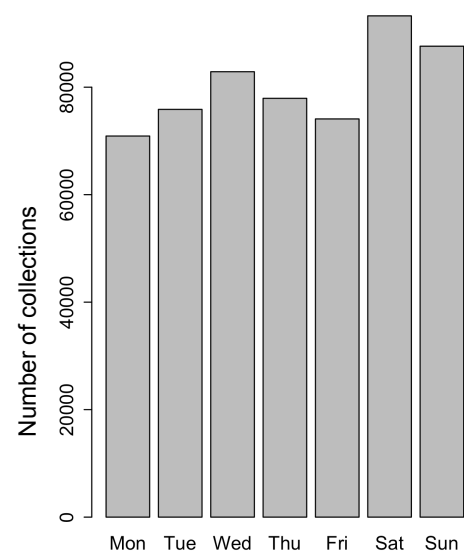
**(b)**



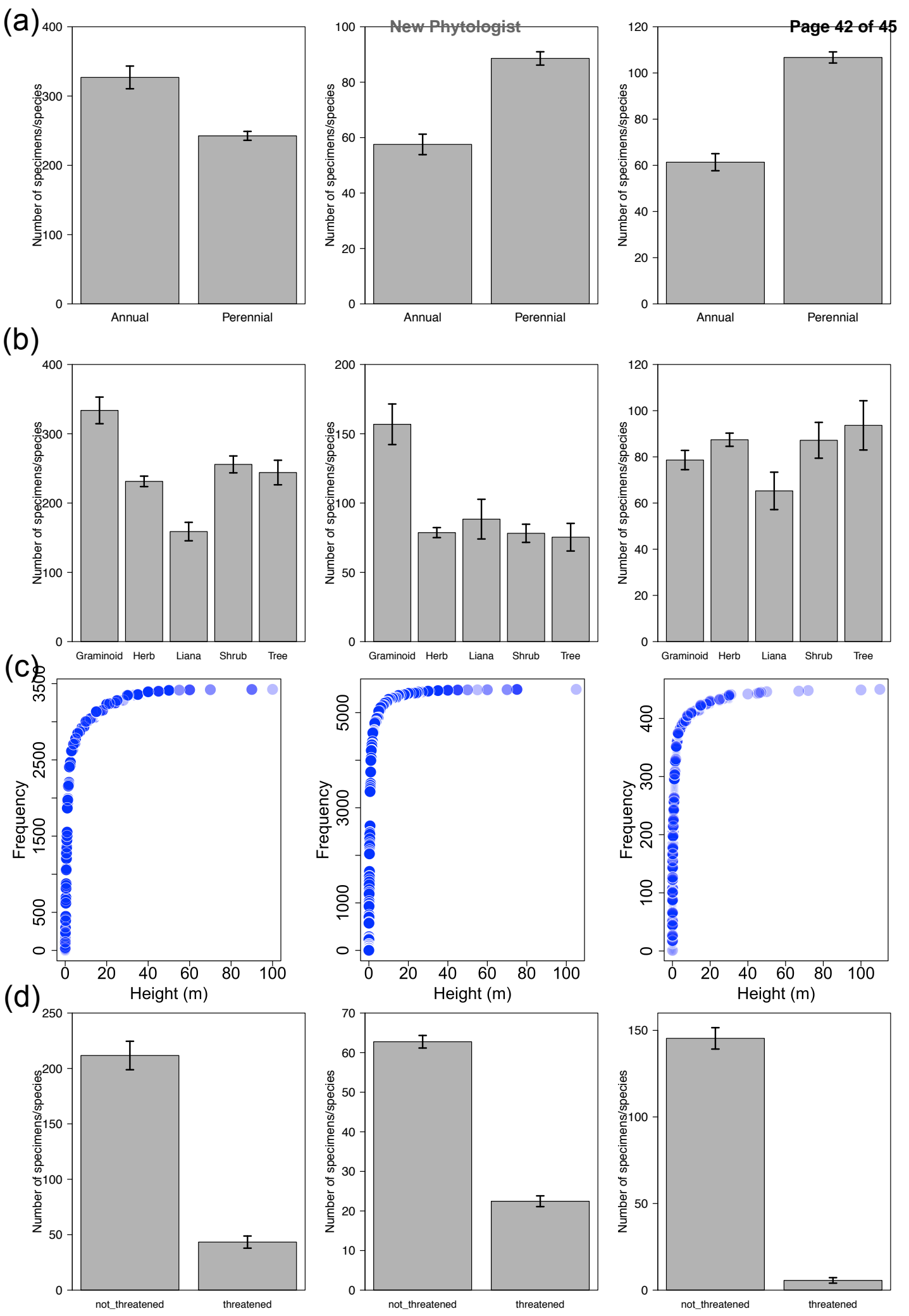
Australia  
(1664 to 2016)



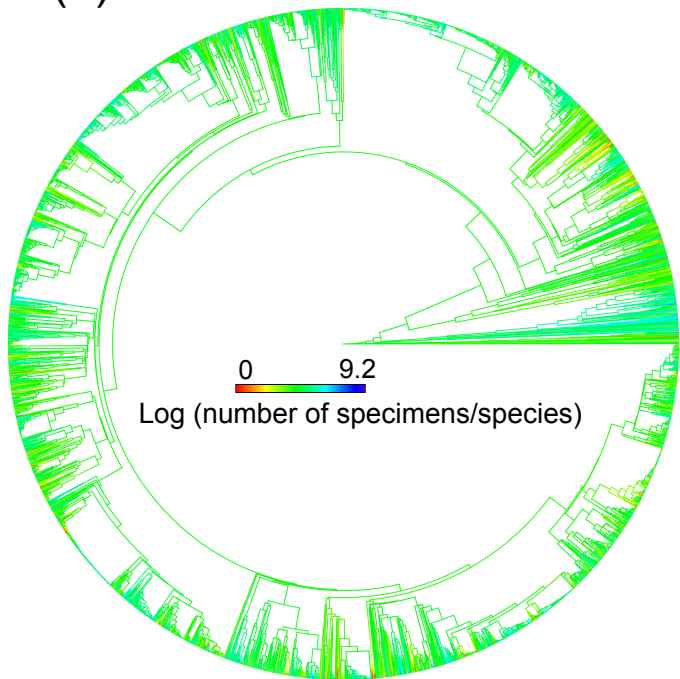
South Africa  
(1656 to 2016)



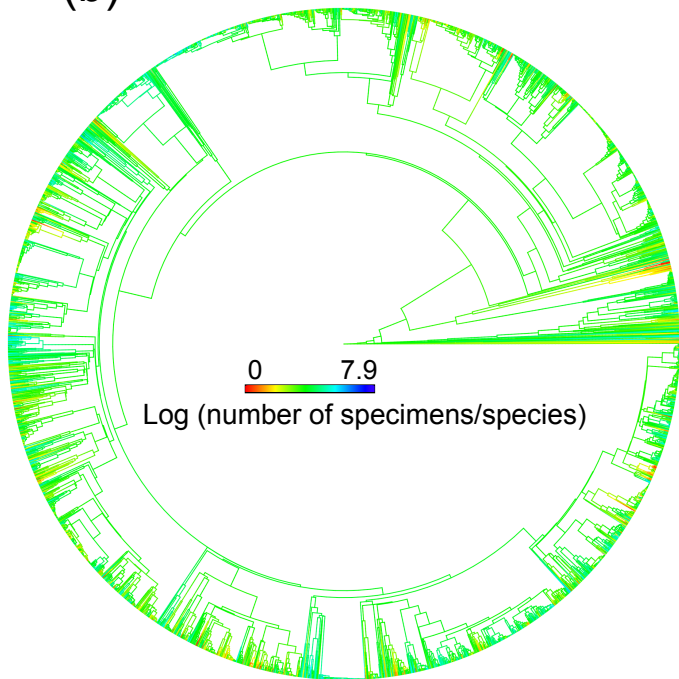
New England  
(1687 to 2016)



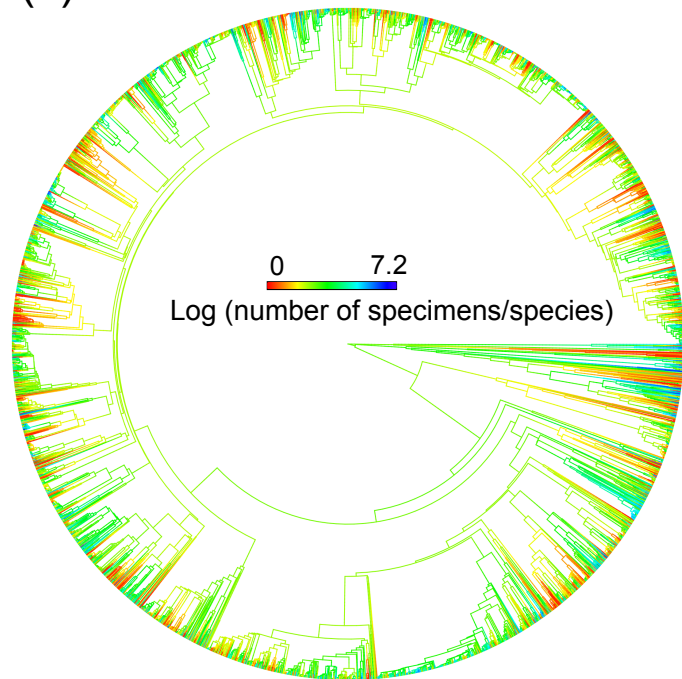
(a)



(b)



(c)



Asteraceae  
 $C_{\text{mean}} = 0.19^{**}$   
 $\lambda = 0.22^{**}$   
 $K = 0.02^{\text{NS}}$

Ericaceae  
 $C_{\text{mean}} = 0.06^{\text{NS}}$   
 $\lambda = 0.16^{\text{NS}}$   
 $K = 0.051^*$

Poaceae  
 $C_{\text{mean}} = 0.06^*$   
 $\lambda = 0.046^*$   
 $K = 0.03^{\text{NS}}$

Brassicaceae  
 $C_{\text{mean}} = 0.02^{\text{NS}}$   
 $\lambda = 0.00007^{\text{NS}}$   
 $K = 0.037^{\text{NS}}$

Fabaceae  
 $C_{\text{mean}} = 0.02^{\text{NS}}$   
 $\lambda = 0.09^{\text{NS}}$   
 $K = 0.013^{\text{NS}}$

Ranunculaceae  
 $C_{\text{mean}} = 0.009^{\text{NS}}$   
 $\lambda = 0.00007^{\text{NS}}$   
 $K = 0.03^{\text{NS}}$

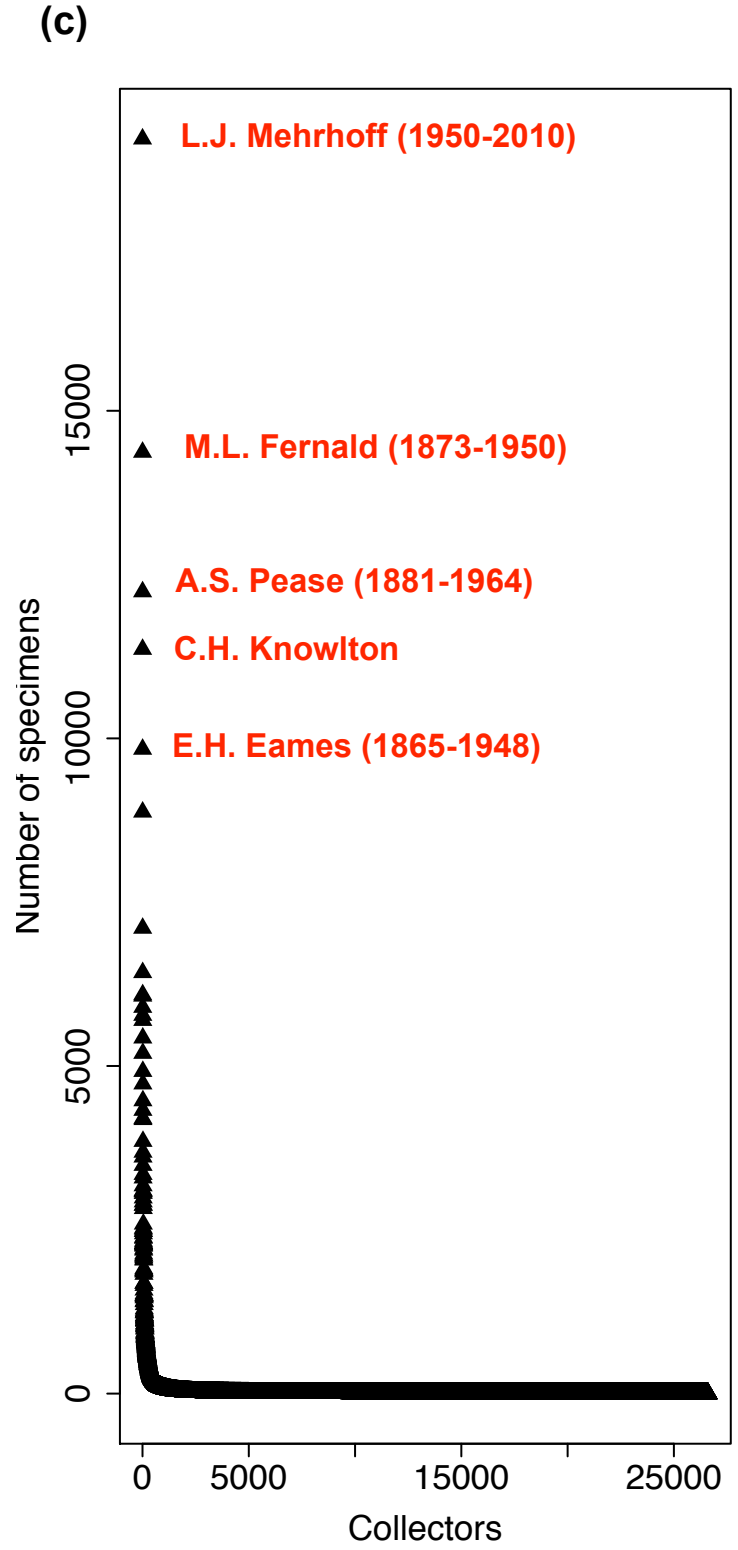
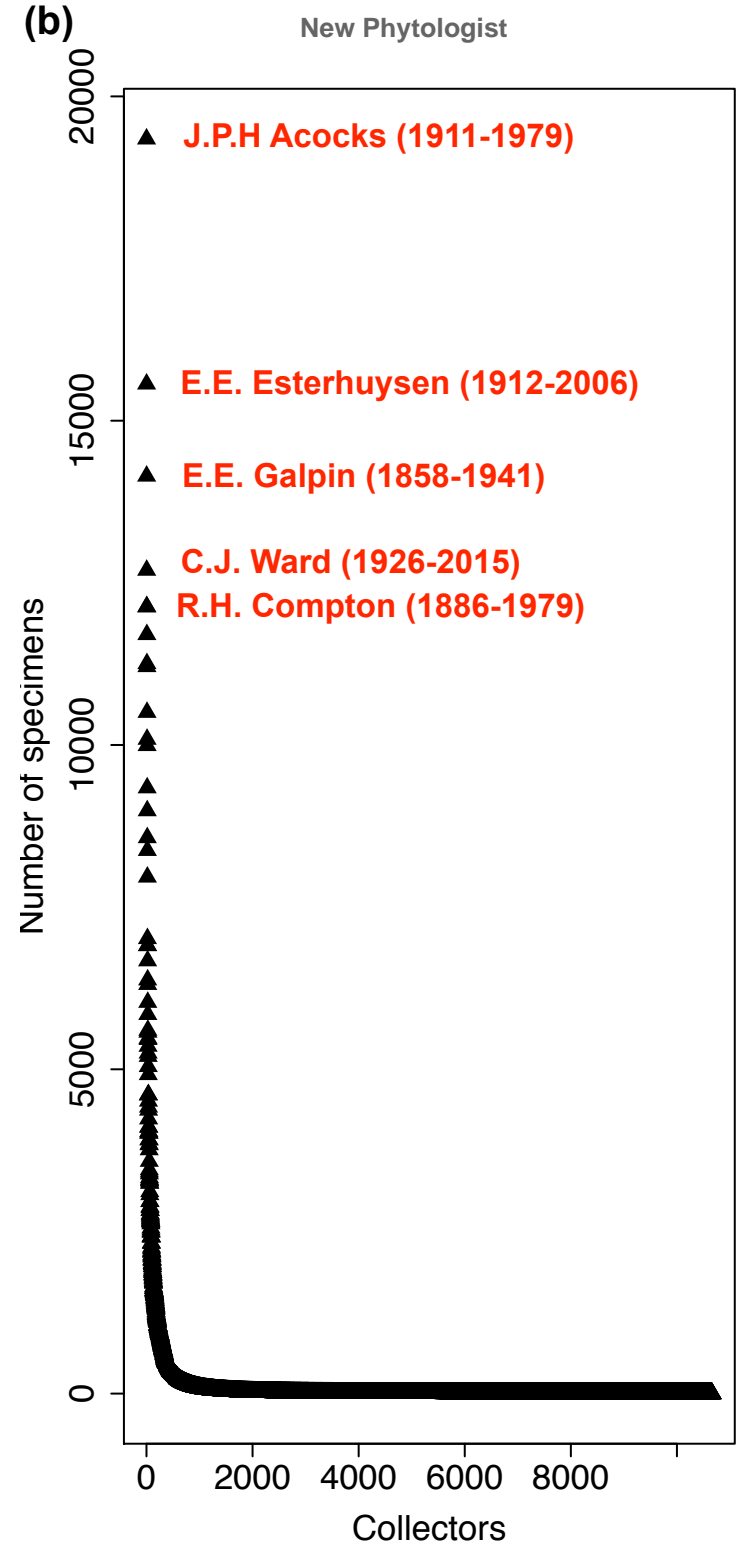
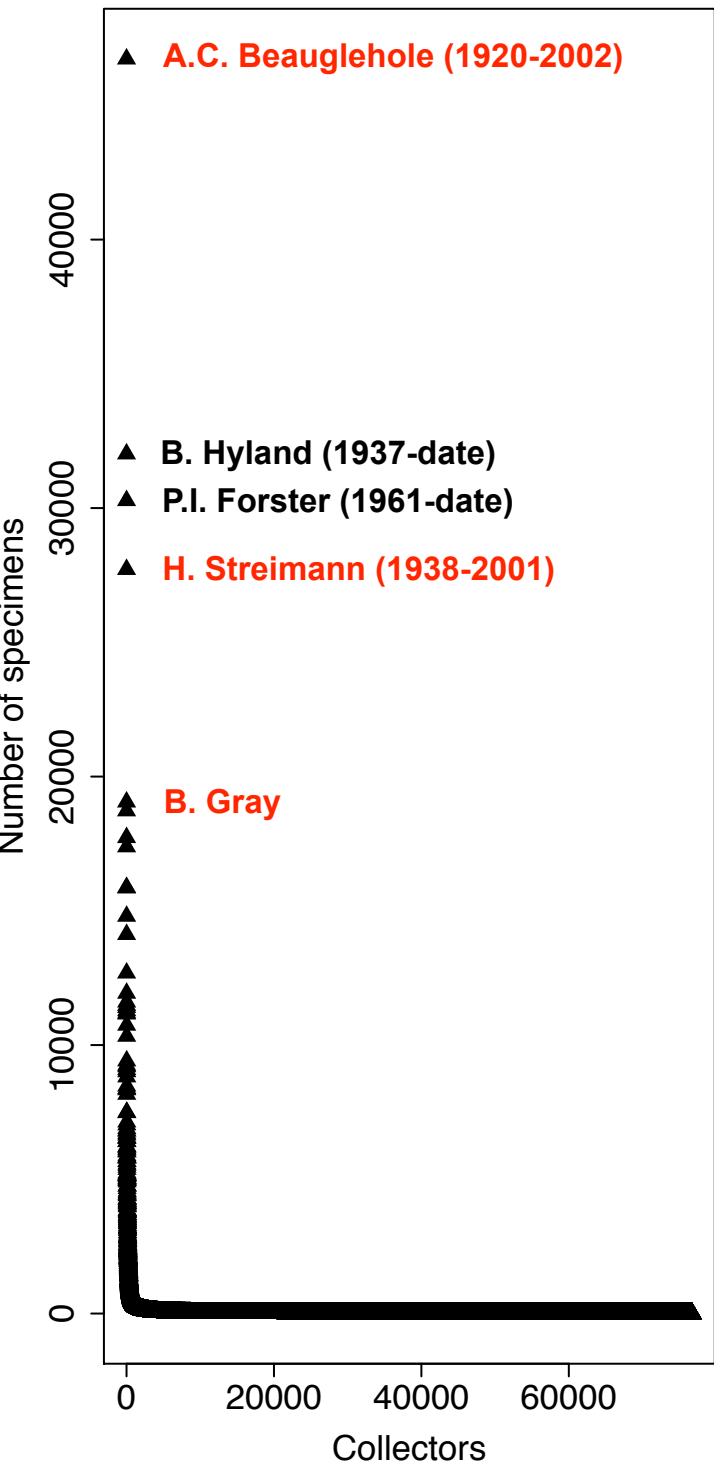
Cyperaceae  
 $C_{\text{mean}} = 0.07^*$   
 $\lambda = 0.035^{\text{NS}}$   
 $K = 0.0077^{\text{NS}}$

Lamiaceae  
 $C_{\text{mean}} = 0.07^{\text{NS}}$   
 $\lambda = 0.21^{\text{NS}}$   
 $K = 0.093^{\text{NS}}$

Rosaceae  
 $C_{\text{mean}} = 0.07^*$   
 $\lambda = 0.00007^{\text{NS}}$   
 $K = 0.009^{\text{NS}}$

Low

High



New Phytologist Supporting Information Figs S1 & S2 and Table S1

Article title: **Widespread sampling biases in herbaria revealed from large-scale digitization**

Authors: Barnabas H. Daru, Daniel S. Park, Richard B. Primack, Charles G. Willis, David S. Barrington, Timothy J. S. Whitfeld, Tristram G. Seidler, Patrick W. Sweeney, David R. Foster, Aaron M. Ellison and Charles C. Davis

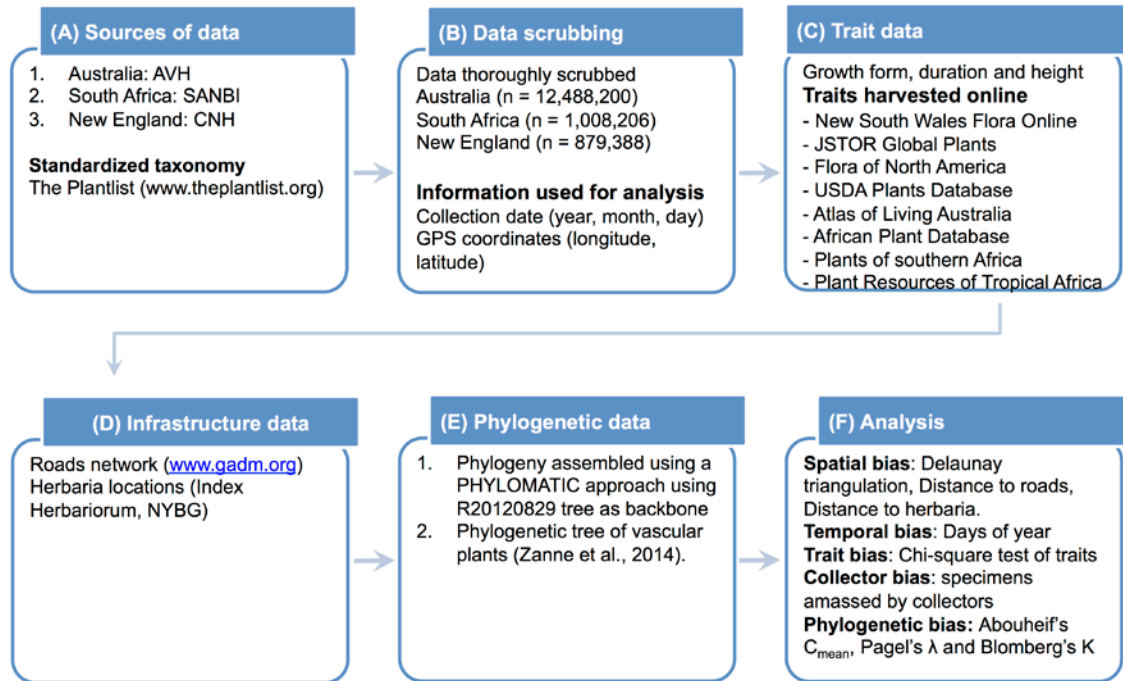
Article acceptance date: 18 September 2017

The following Supporting Information is available for this article:

**Fig. S1** Analytical workflow representing the different steps in the development of this study from data compilation, collation, to statistical analysis.

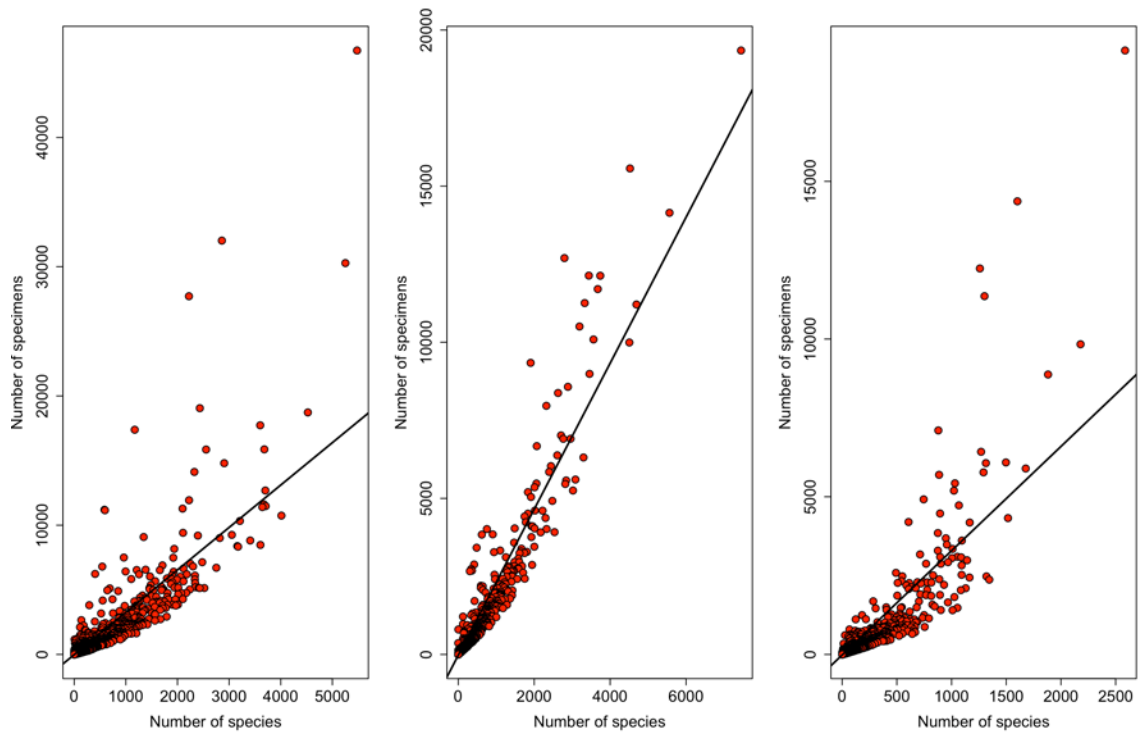
**Fig. S2** Relationships of the number of specimens collected per species with number of species collected in each flora for Australia, South Africa, and New England of the USA.

**Table S1** Results of the tests of phylogenetic signal in the number of specimens collected per species using three methods for nine exemplar clades in New England.



**Fig. S1** Analytical workflow representing the different steps in the development of this study from data compilation, collation, to statistical analysis.





**Fig. S2** Relationships of the number of specimens collected per species with number of species collected in each flora for Australia (left), South Africa (middle), and New England of the USA (right).

**Table S1** Results of the tests of phylogenetic signal in the number of specimens collected per species using three methods (Abouheif's  $C_{\text{mean}}$ , Blomberg's K and Pagel's  $\lambda$ ) for nine exemplar clades in New England: Asteraceae, Brassicaceae, Cyperaceae, Ericaceae, Fabaceae, Lamiaceae, Poaceae, Ranunculaceae, and Rosaceae. Phylogenetic data is derived from Phylomatic (Webb & Donoghue 2005). All tests are based on 1000 randomizations. \*\* $P < 0.001$ ; \* $P < 0.01$ ; NS,  $P > 0.05$

	Asteraceae (n = 593 species)	Brassicaceae (n = 146)	Cyperaceae (n = 518)	Ericaceae (n = 158)	Fabaceae (n = 255 species)	Lamiaceae (n = 146)	Poaceae (n = 565 species)	Ranunculaceae (n = 169)	Rosaceae (n = 346 species)
Abouheif's $C_{\text{mean}}$	0.11**	-0.055 <sup>NS</sup>	0.016 <sup>NS</sup>	0.15*	0.028 <sup>NS</sup>	-0.014 <sup>NS</sup>	0.0026 <sup>NS</sup>	-0.04 <sup>NS</sup>	0.098*
Blomberg's K	0.11 <sup>NS</sup>	0.57 <sup>NS</sup>	0.55 <sup>NS</sup>	0.42 <sup>NS</sup>	0.092 <sup>NS</sup>	0.56 <sup>NS</sup>	0.20*	0.093 <sup>NS</sup>	0.17 <sup>NS</sup>
Pagel's lambda	0.14**	0.00006 <sup>NS</sup>	0.053*	0.27 <sup>NS</sup>	0.23*	0.02 <sup>NS</sup>	0.00008 <sup>NS</sup>	0.00007 <sup>NS</sup>	0.19 <sup>NS</sup>

References to Supporting Information

**Webb CO, Donoghue MJ. 2005.** Phylomatic: tree assembly for applied phylogenetics.

*Molecular Ecology Notes* **5**: 181–183.

**Zanne AE, Tank DC, Cornwell WK, Eastman JM, Smith SA, FitzJohn RG,**

**McGlenn DJ, O'Meara BC, Moles AT, Reich PB, et al. 2014.** Three keys to the radiation of angiosperms into freezing environments. *Nature* **506**: 89–92.