

1 **Identifying the collector practices that shape spatial, temporal, and taxonomic bias in**
2 **herbaria**

3

4 Ryan J. Schmidt^{1*}, Kristen E. Saban¹, Lena Struwe^{2,3}, Charles C. Davis¹

5

6 ¹ Department of Organismic and Evolutionary Biology, Harvard University Herbaria, Harvard
7 University, 22 Divinity Avenue, Cambridge, MA 02138, USA

8 ² Department of Ecology, Evolution & Natural Resources, Rutgers, The State University of New
9 Jersey, 14 College Farm Road, New Brunswick, NJ 08901-8551, USA

10 ³ Department of Plant Biology, Rutgers, The State University of New Jersey, 59 Dudley Road,
11 New Brunswick, NJ 08901-8551, USA

12 * Corresponding author: ryanschmidt@g.harvard.edu

13

14 RJS: 0000-0002-4907-2270

15 KES: 0009-0009-5728-4001

16 LS: 0000-0001-6074-5758

17 CCD: 0000-0001-8747-1101

18

19 **Word Count: 6014**

20 Introduction: 900

21 Materials & Methods: 1301

22 Results: 1962

23 Discussion: 1851

24

25 **Figures:**

26 Manuscript: 8 figures. We would prefer all figures to be in color, however, figures 1 and 6 may
27 be printed in black and white if necessary.

28

29 **Supporting Information: 4 tables and 1 figure.**

30 **Summary**

- 31 • Natural history collections (NHCs) are essential for studying biodiversity. While spatial,
32 temporal, and taxonomic biases in NHCs affect analyses, the influence of collector
33 practices on biases remains largely unexplored.
- 34 • We utilized one million digitized specimens collected in the northeastern United States
35 from 237 herbaria and analyzed contributions from ~10,000 collectors. We investigated
36 (a) similarities and differences between more- and less-prolific collectors, and (b) how
37 these practices influence spatial, temporal, and taxonomic biases.
- 38 • We identified six common collector practices, or collection norms: collectors generally
39 collected (a) different species, (b) from multiple locations, (c) from sites sampled by
40 others, (d) during the principal growing season, (e) species identifiable outside peak
41 collecting months, and (f) species from species-poor families and genera. Some norms
42 changed over decades, with different taxa favored during different periods. Collection
43 norms have increased taxonomic coverage in NHCs, however, collectors typically
44 avoided large, taxonomically-complex groups, causing their underrepresentation in
45 NHCs. Less-prolific collectors greatly enhanced coverage by collecting during more
46 months and from less-sampled locations.
- 47 • We assert that overall collection biases are shaped by shared predictable collection
48 norms rather than random practices of individual collectors. Predictable biases offer an
49 opportunity to more effectively address biases in future biodiversity models.

50
51 **Keywords**

52 herbaria; natural history collections; history of science; collection norms; biodiversity;
53 digitization; biodiversity modelling

54
55 **Introduction**

56 Discovering and describing global patterns of species diversity and distribution remains a
57 fundamental priority for biodiversity scientists (CBD, 2022). Although recent advances in
58 biodiversity modeling have greatly improved our understanding of these factors, the vouchered
59 specimens and observational data underlying these models are known to exhibit significant
60 spatial, temporal, and taxonomic biases that remain largely unaccounted for (Meyer *et al.*, 2016;
61 Daru *et al.*, 2018).

62

63 Herbaria and other natural history collections (NHCs) are invaluable resources for
64 understanding global biodiversity (Funk, 2003; Johnson *et al.*, 2023; Davis, 2023, 2024; Marín-
65 Rodulfo *et al.* 2024). The extensive sampling of NHCs over time, space, and taxa complement
66 long-term monitoring programs such as the Atlas of the British Flora (Perring & Walters, 1962;
67 Preston, 2013) and the USDA’s Forest Inventory and Analysis (Rudis, 2003; FIA, 2023), which
68 have provided important insights into species distributions but are limited across these key axes
69 in important ways. Although biodiversity is not randomly distributed, to best represent
70 biodiversity NHCs would ideally aim to represent as close to an unbiased sample of global
71 biodiversity across time, space, and taxa as possible. Understanding how NHCs diverge from
72 these ideals allows us to better account for biases in our biodiversity models and discern what
73 questions we can address using these collections. Ultimately, understanding collection biases
74 will help guide the application and development of statistical tools to correct for biases, develop
75 better priorities for future collecting efforts, and help us achieve more comprehensive and
76 accurate models of global biodiversity.

77
78 Comprehensive digitization of natural history specimens from large geographic/floristic regions
79 has revealed key spatial, temporal and taxonomic biases in NHCs (Meyer *et al.*, 2016; Daru *et*
80 *al.*, 2018; Eckert *et al.*, 2024). These overall biases in NHCs are a consequence of the spatial,
81 temporal, and taxonomic collection practices of each collector—what we call collector practices.
82 Previous studies have highlighted the connection between collector practices and overall bias in
83 collections, documenting that a small number of mega-collectors have made disproportionately
84 large contributions to species discovery (Bebber *et al.* 2012) and to specimen collections in
85 NHCs (Daru *et al.* 2018). The disproportionately large impact of these mega-collectors raises an
86 important but unanswered question: have highly prolific collectors also contributed
87 disproportionately to the biases documented in these collections? To date, there have been no
88 efforts to investigate how the collector practices of all collectors in a region have contributed to
89 overall bias in NHCs. Moreover, there have been no large-scale efforts to understand the impact
90 that less-prolific collectors have had on the spatial, temporal, and taxonomic coverage in
91 collections.

92
93 Here, we examine the impact of collector practices on novel and previously documented biases
94 in NHCs (Meyer *et al.*, 2016; Daru *et al.*, 2018; Kozlov *et al.*, 2021; Eckert *et al.*, 2024). As a test
95 case for our investigation, we leverage the nearly completely digitized metaherbarium that
96 extensively documents the flora of the northeastern United States (Schorn *et al.*, 2016;

97 Sweeney *et al.*, 2018; Hedrick *et al.*, 2020). Specifically, we use all digitized herbarium
98 specimens of land plants (i.e., bryophytes and vascular plants) collected in the northeastern
99 United States from the earliest digitized record to the present (i.e., 1781–2024). We reconstruct
100 the contributions of collectors to investigate how overall bias in NHCs is impacted by the
101 similarities and differences in collection practices of different collectors. We assess the
102 relationship between these collection practices and the number of collections by each collector
103 on a continuous scale with more- and less-prolific collectors representing opposite ends of this
104 continuum. Mega-collectors—who have contributed a disproportionately large amount of
105 specimens (*sensu* Daru *et al.*, 2018)—represent the uppermost extreme of this spectrum. We
106 also investigate how what we term *collection norms*—the collector practices shared by all
107 collectors—have influenced overall biases in NHCs. Such synthetic investigations further
108 demonstrate the growing utility of digitized specimens within the framework of the extended
109 specimen (Webster, 2017; Lendemer *et al.*, 2020), facilitating proper attribution for the
110 thousands of hidden heroes that have made meaningful but previously unrecognized
111 contributions to NHCs (Groom *et al.*, 2022) and enabling ongoing efforts to better model
112 biodiversity in an era of rapid ecological change.

113

114 **Materials and Methods**

115 ***Data collection & data cleaning***

116 We downloaded 2,365,287 records representing all digitized herbarium specimens of land
117 plants from the northeastern United States (i.e., Connecticut, Maine, Massachusetts, New
118 Hampshire, New Jersey, New York, Pennsylvania, Rhode Island, and Vermont; hereafter the
119 Northeast) from GBIF (GBIF.org, 2024). These specimens are housed in 237 herbaria around
120 the world (Table S1). We then filtered this dataset to remove the 548,895 records without a
121 transcribed date, collector, locality, or species-level identification. This filtering left us with
122 1,816,392 analyzable records.

123

124 ***Georeferencing***

125 About half of the cleaned records (920,633 records) contained transcribed coordinates. We
126 batch-georeferenced an additional 401,450 specimens to municipal centroid points (CT DEEP,
127 2023; PennDOT, 2024) and removed all records that could not be georeferenced to a specific
128 municipality (503,563 records removed). Although this method of georeferencing does not
129 capture fine-scale differences in collection localities (Park & Davis, 2017), it is consistent with
130 the precision for many herbarium georeferencing initiatives in the northeastern US (e.g., Mancini

131 *et al.*, 2019) and suitable for analyses on these large spatial scales. We removed records with
132 coordinates outside of the northeastern US (United States Census Bureau, 2024) using the
133 `st_intersection()` function from the *sf* package in R 4.4.1 (Pebesma 2018; Pebesma &
134 Bivand 2023; 10,254 records removed). This resulted in a total of 1,311,829 georeferenced
135 records.

136

137 **Collector disambiguation**

138 Due to institutional differences in transcription practices, incorrect transcriptions, and
139 orthographic variations in collector names, assigning different text strings (i.e., recordedBy
140 strings in DarwinCore; hereafter “collector strings”) to a single collector can be difficult and time
141 consuming for large datasets (Groom *et al.*, 2022). Thanks to the large-scale availability of
142 digitized historical and genealogical records (e.g., Ancestry.com, MyHeritage.com, and
143 Newspapers.com) and recent initiatives by historians of science to identify and disambiguate the
144 names of people who collected natural history specimens (e.g., Bionomia; Shorthouse, 2024;
145 Weeks *et al.*, 2024), we are for the first time able to identify and reconstruct the oeuvres of all
146 contributors to a regional flora.

147

148 To disambiguate collector strings, we extracted the first collector in each collector string,
149 separating what we consider the principal collector (henceforth referred to as the collector) from
150 any associated collectors. Although associated collectors are crucial parts of any collection
151 team and deserve proper credit for their efforts, we focused our analysis on principal collectors
152 in this initial study. Our rationale is that the principal collector is usually responsible for recording
153 field notes and is likely to take on the major role of depositing the specimens in an herbarium
154 collection. We then separated the collector strings into words using the `unnest_tokens()`
155 function from `tidytext` (Silge & Robinson, 2016) and concatenated these words in
156 alphabetical order to standardize different transcriptions of the same text (e.g., “C. F. Parker”, “C
157 F Parker”, and “Parker, C. F.” would all become “c,f,parker”). We then merged all records with
158 identical concatenated strings and manually validated each cluster—merging records with
159 different concatenated strings that represent the same collector—to ensure that each cluster
160 represented a single collector. We used biographical information from historical and
161 genealogical databases (e.g., Ancestry.com and Newspapers.com) and databases of natural
162 history collectors (i.e., Bionomia and Harvard Index of Botanists; Shorthouse, 2024; Harvard
163 University Herbaria, 2024) to reconstruct the oeuvres of collectors that collected under multiple
164 names, including their spouses’ names. For instance, we identified “Mrs. C. S. Phelps” as Ora

165 Almira Phelps (née Parker) who collected under the names Mrs. Charles Sheppard Phelps,
166 Orra A. Phelps, Mrs. O. P. Phelps, and Orra Parker Phelps.

167

168 We excluded any collector strings that were ambiguous either because of obvious transcription
169 errors that could not be verified with a digital image of the specimen or had limited information.

170 To ensure that we were not conflating multiple collectors, we excluded records with only initials

171 (e.g., C.A.B.), only a surname (e.g., Boice), or only the initial of the first name and the surname

172 (e.g., C. Boice; 233,321 records removed, 1,078,508 records remaining). We then removed

173 duplicate specimens (i.e., specimens collected by the same collector with the same specimen

174 number in DarwinCore's `recordNumber` field) so that each collection event is represented by a

175 single specimen (89,251 records removed). This resulted in our final dataset of 989,257

176 specimens (Table S2).

177

178 ***Temporal Bias***

179 To investigate temporal trends in botanical collections, we calculated the number of specimens,

180 distinct species, sampling localities, and active collectors for each year during 1781–2024. We

181 investigated seasonal variations in collection intensity by comparing the number of specimens

182 collected in each month and analyzed how this distribution changed with respect to the oeuvre

183 size of the collector who gathered the specimen.

184

185 ***Spatial Bias***

186 We quantified spatial bias by gridding the georeferenced specimens into 10-km grid squares

187 (hereafter localities) to help mitigate the effects of batch georeferencing and create equal-area

188 polygons for comparison (Franklin & Miller, 2009; Schmidt *et al.*, 2023). We calculated the

189 revisitation proportion for each collector as the number of specimens per unique collecting

190 locality. We also calculated the average oeuvre size of collectors active in each locality,

191 weighted by the number of collections of each collector (higher values indicate more activity by

192 highly prolific collectors) to investigate the geographical bias of more- versus less-prolific

193 collectors.

194

195 To understand how collectors of different sizes contributed to overall spatial sampling, we found

196 the number of unique grids sampled for different subsets of the data. To determine if more- or

197 less-prolific collectors expand overall spatial coverage, we arranged specimens by decreasing

198 and increasing oeuvre size, respectively and found the number of unique grids sampled for

199 increasingly larger subsets of the data in 10,000 specimen increments (i.e., after arranging by
200 oeuvre size, we extracted the first 10,000 specimens, first 20,000 specimens, 30,000
201 specimens, etc.). We assessed how spatial bias from collectors with different oeuvre sizes
202 differs from two different null models: we randomly ordered specimens from our dataset to
203 determine if collections by more- or less-prolific collectors are more spatially clustered than the
204 overall specimens (randomized specimens); and simulated a new dataset by randomly sampling
205 from all localities in the northeastern US to determine how collections differ from spatially
206 random collections (simulated random sampling).

207

208 ***Taxonomic Bias***

209 To determine the relative representation of different taxa in herbarium collections, we calculated
210 collection depth as the average number of specimens per species in a given taxon in the
211 northeastern US (i.e., total specimens/unique species for each genus and family). We evaluate
212 taxon size on a continuous scale, whereby taxa with fewer species in the northeastern US are
213 considered smaller and those with more species are considered larger. Taxa with higher
214 collection depths were considered better represented in herbaria.

215

216 To assess how frequently collectors collect a species that they have already collected, we
217 calculated the proportion of species re-collected by each collector (i.e., total specimens/unique
218 species for each collector). Collectors who collected many specimens of the same species
219 would have a high re-collection proportion while those that collected only one specimen of each
220 species would have a re-collection proportion of one.

221

222 To investigate whether some taxa (i.e., species, genera, and families) were favored by
223 collectors over other taxa, we plotted the number of collections per taxon against the number of
224 collectors who collected each taxon. We fit a generalized additive model (GAM) to these points
225 to estimate how many collectors we expected to have collected each taxon based on the total
226 number of specimens of that taxon. Taxa that fell above this GAM curve were collected by more
227 people than expected (hereafter, favored taxa) and taxa that fell below the curve were collected
228 by fewer people than expected (hereafter, commonly avoided taxa).

229

230 **Results**

231

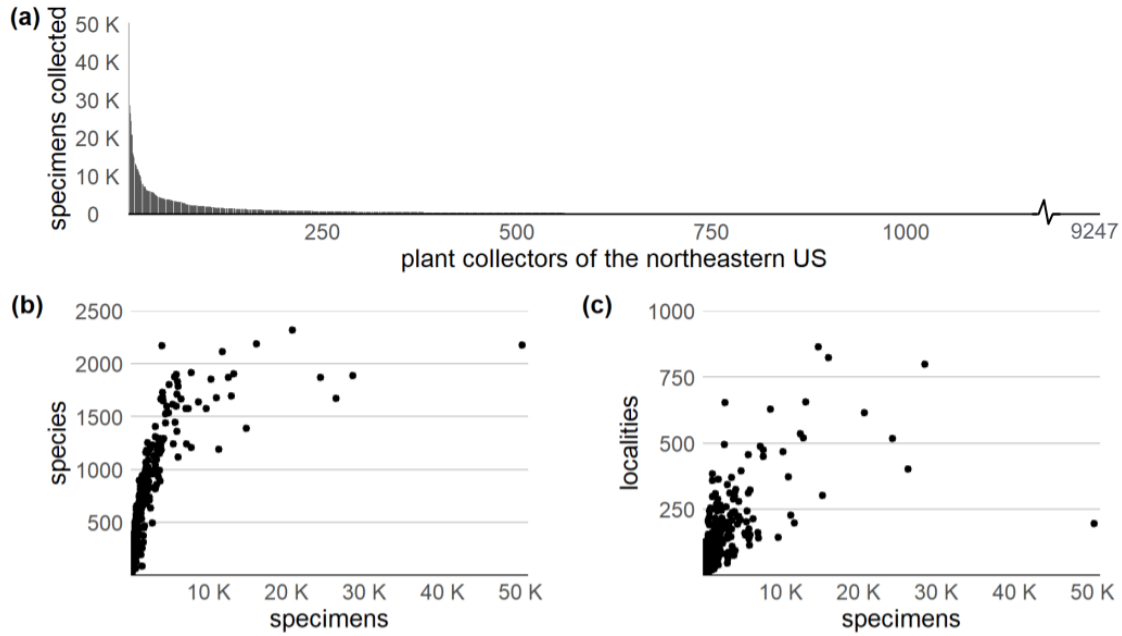
232 ***Collectors***

233 We identified 9247 collectors who collected plant specimens in the northeastern US (Fig. 1a;
234 Table S3). This is no doubt an underestimate of the total number of people who have
235 contributed to collections in the region since many collectors were excluded from our analysis
236 due to incomplete or ambiguous collector names and insufficient locality information (45% of
237 analyzable specimens removed) and those whose specimens have yet to be digitized. There
238 was a large variation in the number of specimens that each collector collected. The vast majority
239 (more than 90%; 8385 people) collected fewer than 100 specimens. Only 1.8% of collectors
240 (171 people) collected more than 1000 specimens (contributing 71% of the total number of
241 collections). The most prolific collector in our dataset was Raymond L. Schaeffer, Jr., who
242 collected 50,287 specimens (Fig. 1b). Half of all specimens from the northeastern US were
243 collected by only 57 collectors (0.6% of collectors). Most collectors (70%; 6,549 people)
244 collected fewer than ten specimens (contributing 1.5% of collections).

245

246 People who collected less than 1000 specimens tended to collect only one specimen of each
247 species (Fig. 1c) and about ten specimens per locality (Fig. 1d). For collectors who collected
248 more than 1000 specimens, they tended to collect only one specimen for each species for the
249 first 1000 specimens they collected. After collecting about 1000 specimens, they collected
250 multiple specimens of the same species, and the number of species they collect plateaus near
251 2000 species. E. H. Eames collected the most plant species of any collector in our dataset
252 (2574 species, both vascular and nonvascular; Whelan, 1948). Most people collected either
253 vascular plants (85%) or non-vascular plants (7%), with only 8% collecting both types.

254



255

256

257

258

259

260

261

262

263

Temporal Bias

264

265

266

267

268

269

270

271

272

273

274

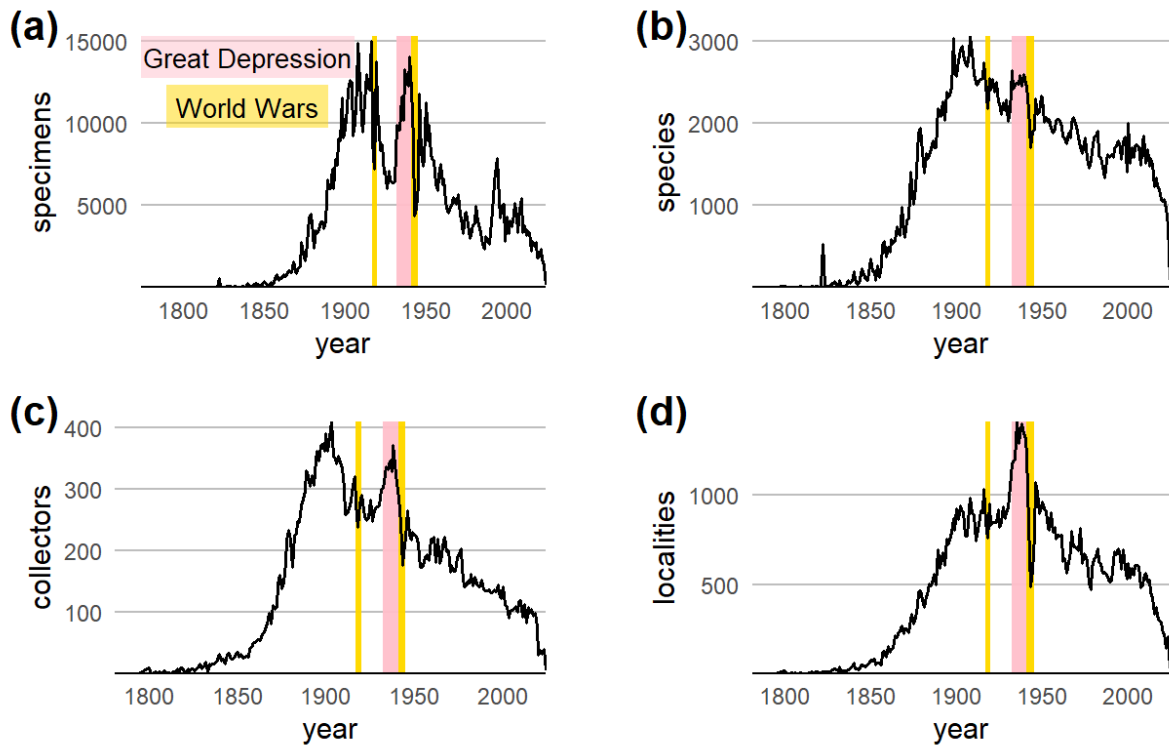
275

Figure 1. We identified 9247 people who collected herbarium specimens in the northeastern US. The bar plot shows (a) the total number of unique specimens for each plant collector in the northeastern US. The scatter plots show the relationship between the number of specimens each person collected and (b) the number of species they collected and (c) the number of localities in which they collected.

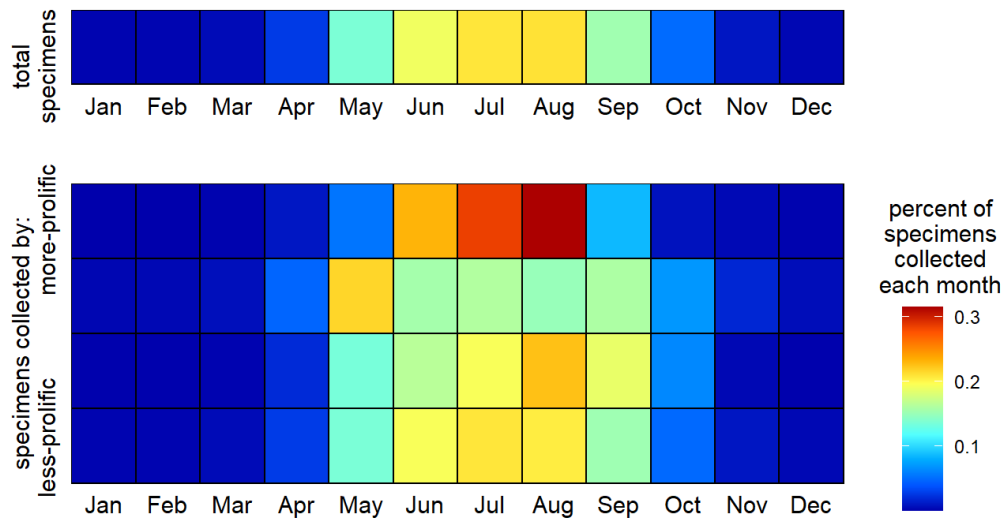
The number of collectors active in a given year has varied substantially through time with peaks during 1880–1916 and again during 1932–1941 (Fig. 2a–d). The number of active collectors is strongly correlated with the number of specimens collected in a given year (cross-correlation value of 0.90, $p < 0.001$), species (0.94, $p < 0.001$), and localities (0.90, $p < 0.001$). The number of specimens (Fig. 2a) and the number of species collected (Fig. 2b) in a given year also peaked during 1880–1916 and 1935–1941 whereas the number of sampling localities peaked only from 1935–1941 (Fig. 2d). All metrics have declined since 1950.

About 90% of specimens from the northeastern US were collected during spring and summer (i.e., May to September)—the main growing season in northern temperate zones—with relatively few specimens collected during off-peak months (i.e., from October through April; Fig. 3). The highest proportion of collections by less-prolific collectors were also during May–

276 September. However, collections by more-prolific collectors had a much narrower temporal
277 distribution with collections almost exclusively from June, July, and August.
278



279
280 **Figure 2.** The line plots show the annual variation in (a) the number of specimens
281 collected, (b) the number of species collected, (c) the number of active collectors, and
282 (d) the number of localities in which specimens were collected from 1781–2024. The
283 yellow bars indicate the years when the US was involved in World Wars I and II (1917–
284 1920 and 1941–1946, respectively) and the pink bars represent the federal
285 environmental projects during the Great Depression (1929–1939).
286

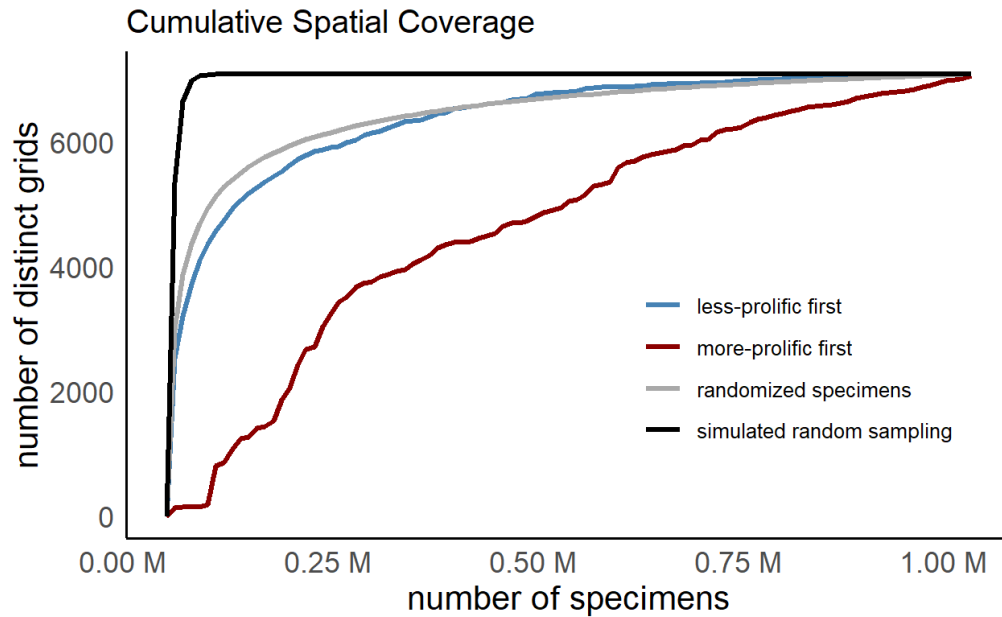


287
 288 **Figure 3.** This graph shows the percentage of specimens collected in each month for all
 289 specimens (total specimens) and subdivided into four bins based on oeuvre size (i.e.,
 290 going from the 25% of specimens collected by the least prolific collectors at the bottom
 291 to the 25% of specimens collected by the most prolific collectors at the top).

292
 293 **Spatial Bias**

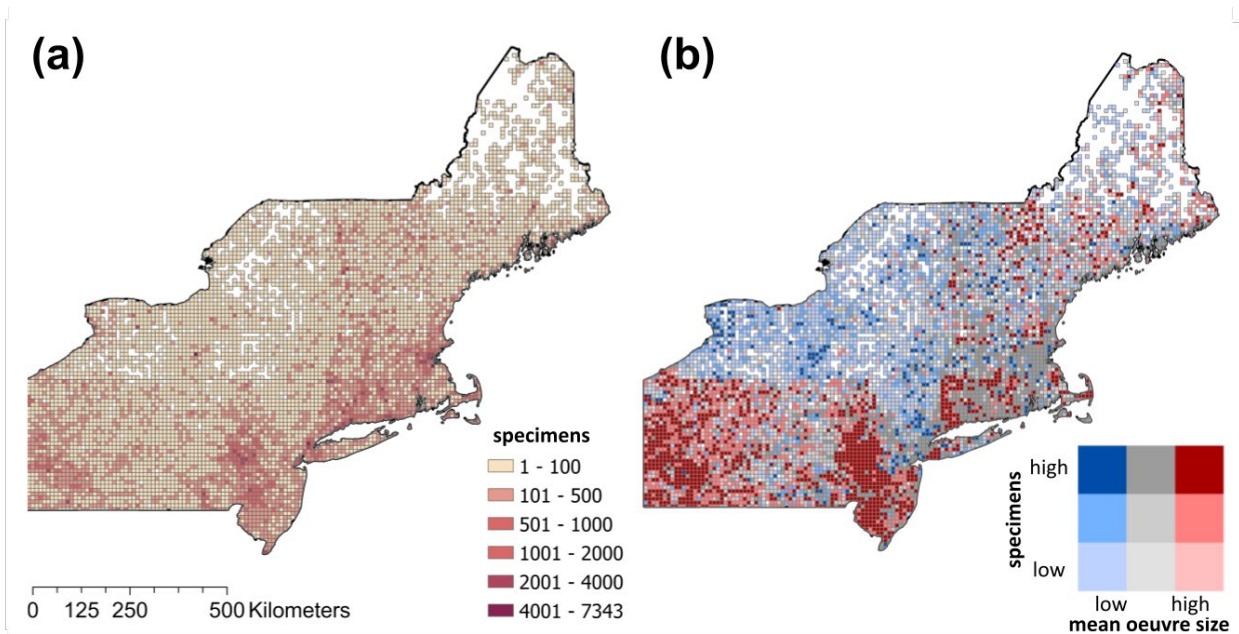
294 The specimens collected by more-prolific collectors were more spatially clustered and had lower
 295 geographic coverage than those collected by less-prolific collectors (Fig. 4). Additionally,
 296 collections by less-prolific collectors included areas not represented by more-prolific collectors,
 297 but more-prolific collectors did not capture areas not represented by less-prolific collectors.

298
 299 Certain spatial clusters that dominate overall specimen clustering in the northeastern US are
 300 driven almost exclusively by collections from more-prolific collectors (Fig. 5). Some of the areas
 301 with the highest collection density are driven by a few, prolific collectors (e.g., the hotspot in
 302 near Allentown, PA is driven primarily by R. L. Shaeffer, Jr.), whereas other areas with high
 303 collection density are driven by many less-prolific collectors (e.g., many of the hotspots in
 304 upstate NY). The overall density of collections and the different drivers of collection intensity
 305 change quickly over some state borders. For example, there are dense collections in PA and
 306 very sparse collections in adjacent NY.



307
 308
 309
 310
 311
 312
 313
 314
 315

Figure 4. Accumulation curves for the cumulative spatial coverage of gridded herbarium specimens based on the oeuvre size for collectors. Specimens were added by decreasing oeuvre size for the red curve (more-prolific collectors added first); increasing oeuvre size for the blue curve (less-prolific collectors added first); and in a random order independent of oeuvre size for the gray curve (randomized specimens, median of 99 permutations). The black curve shows randomly simulated specimens to represent our null model of random spatial sampling in the region (simulated random sampling).



316

317 **Figure 5.** The maps show (a) the density of collections in the northeastern US and (b)
318 the relationship between collection density and areas where collections have been
319 driven primarily by less-prolific collectors (blue; bottom 33% of collectors with the
320 smallest oeuvres), more-prolific collectors (red; top 33% of collector oeuvres), or a mix of
321 collector types (gray; middle 33% of collector oeuvres).

322

323 ***Taxonomic Bias***

324 Smaller genera are more likely to have a greater collection depth than larger genera; the same
325 is the case for smaller families (Fig. 6). Despite the overrepresentation of smaller genera,
326 several of the most frequently collected species are from large genera (e.g., three species of
327 *Carex*; for a list of the hundred most frequently collected species, see Table S4). Ferns were
328 dramatically overrepresented among the most frequently collected species (11 of the top 20
329 collected species were ferns). Within each year, 90% of specimens were collected during May–
330 September but only 46% of species were collected only during these five months. Species that
331 have been collected outside of the peak collection window (i.e., with at least one collection
332 during October–April) are far more likely to be overrepresented in herbaria compared with
333 species that have not been collected outside of peak collection months (Fig. S1). These non-
334 peak species include all but 18 of the 1000 most commonly collected species in the Northeast;
335 11 of these 18 are species of *Carex*. Despite also being collected in off-peak months, the top
336 species have been preferentially collected throughout the year, including during peak months;
337 96% of the top 1000 most collected species remain in the top 1000 when only collections from
338 peak months are considered.

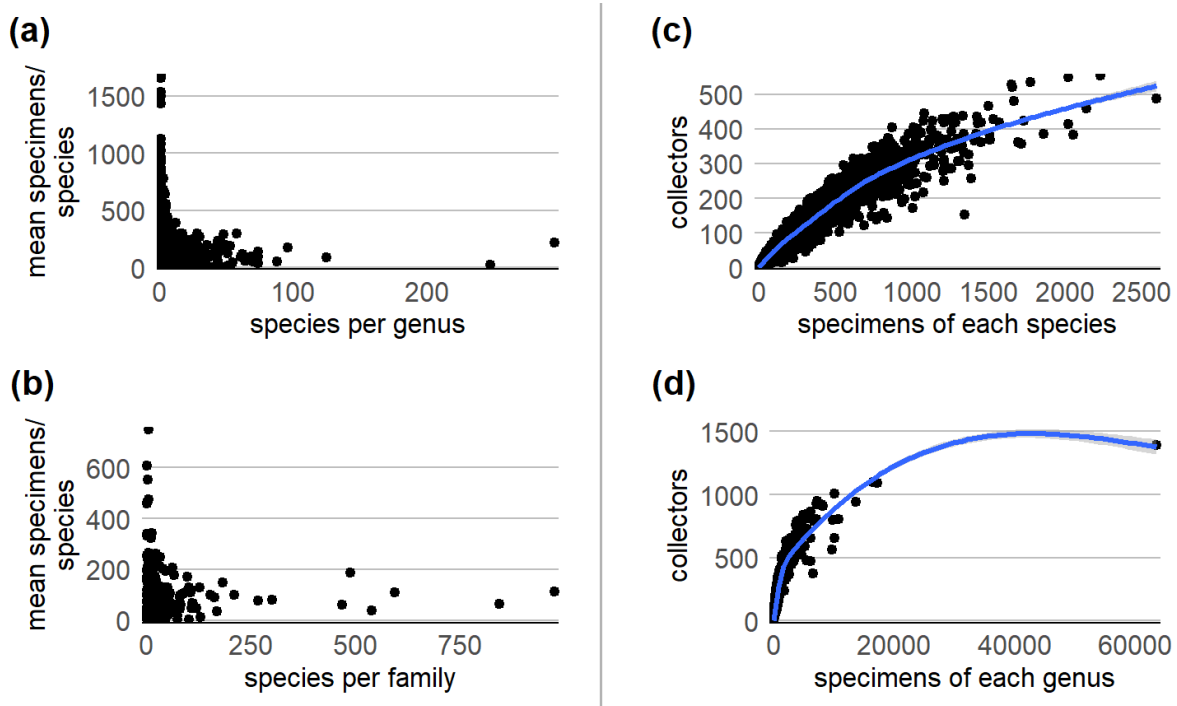
339

340 Some species are overrepresented in collections because they were collected by many people
341 (e.g., *Arisaema triphyllum* (L.) Schott, *Onoclea sensibilis* L., and *Polystichum acrostichoides*
342 (Michx.) Schott), whereas others are overrepresented because they were collected intensively
343 by a few people (e.g., *Sceptridium dissectum* (Spreng.) Lyon, *Scirpus cyperinus* (L.) Kunth, and
344 *Viola sororia* Willd.).

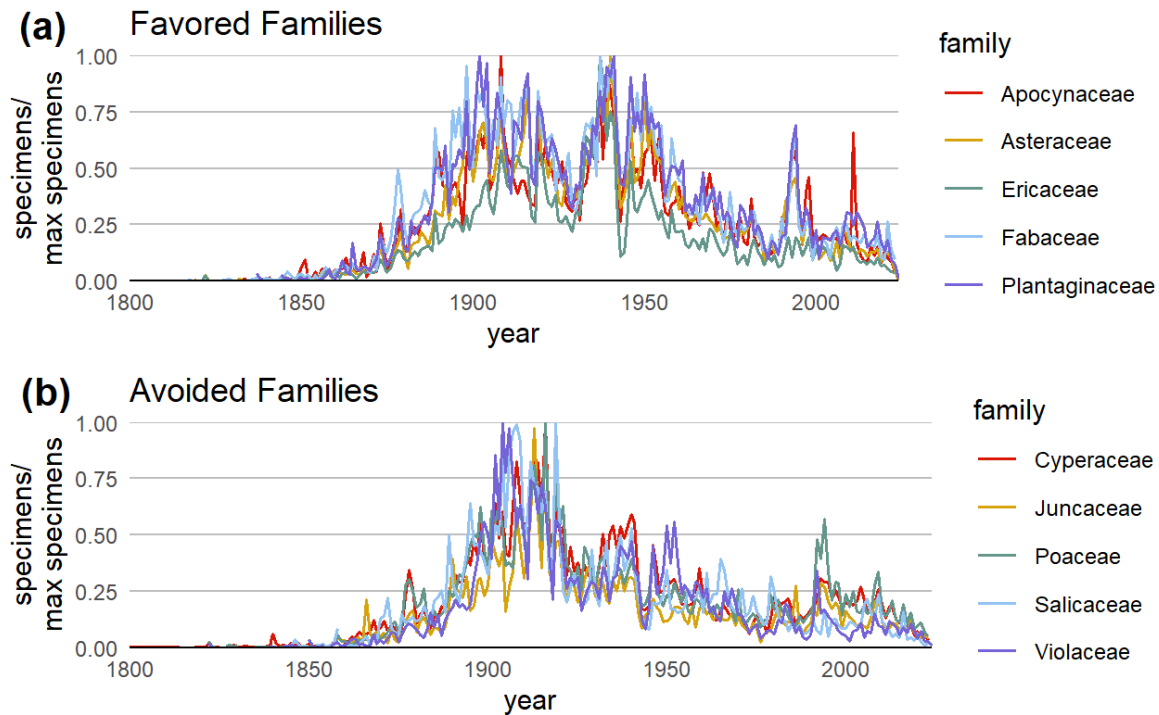
345

346 Some species were collected by far more people than expected from our GAM model (e.g.,
347 *Cypripedium acaule* Aiton and *Solanum dulcamara* L.) whereas *Dichantherium acuminatum*
348 (Sw.) Gould & C.A.Clark was collected by far fewer people than expected. Similarly, some
349 genera were collected by more people than expected from our GAM model (e.g., *Lobelia*,
350 *Lysimachia*, and *Trifolium*), whereas others by fewer people than expected (e.g., *Crataegus*,

351 *Dichanthelium, Potamogeton, Salix, Sphagnum*). Some families were also collected by more
352 people than expected from our model (e.g., Apocynaceae, Asteraceae, Ericaceae, Fabaceae,
353 and Orchidaceae) and others by fewer than expected (e.g., Cyperaceae, Poaceae, Juncaceae,
354 Salicaceae, and Violaceae). Commonly favored families—collected by more people than
355 expected—typically had peaks in annual collections in the 1910s and 1930s, mirroring overall
356 trends in collections through time (Fig. 7). Commonly avoided families—collected by fewer
357 people than expected—typically had only a single peak during the 1910s. Some commonly
358 avoided families (e.g., Potamogetonaceae and Sphagnaceae), had relatively low collections
359 through time and its peaks correspond to specialist collectors rather than overall trends in
360 collections.
361



362
363 **Figure 6.** The plots show the collection depth (average number of specimens per
364 species) for each (a) genus and (b) family. The scatter plots in the right pane (panels c &
365 d) show the relationship between the number of specimens per species and the number
366 of collectors who collected these species of each (c) species and (d) genus.
367



369
 370 **Figure 7.** The annual variation in collection intensity for a subset of families collected by
 371 (a) more people than expected (favored families) and (b) less people than expected
 372 (avoided families). The vertical axes are adjusted to show variation in collection intensity
 373 for each family on the same scale where 1 represents the maximum number of
 374 specimens collected in a given year for each family.

375

376 **Summary of Results**

377 We identified nearly 10,000 collectors who have made important contributions to our
 378 understanding of plant biodiversity in the northeastern United States. We confirmed that a few
 379 mega-collectors contributed a disproportionately large share of these collections. Our analysis
 380 reveals many novel ways in which the collection efforts by thousands of less-prolific collectors
 381 have greatly expanded the temporal, spatial, and taxonomic dimensions of NHCs.

382

383 We assert that overall bias in collections across space, time, and taxa, is strongly impacted by
 384 predictable collection norms that are the result of the shared collector practices of many
 385 collectors rather than by stochastic biases of individual collectors (Fig. 8). The predictability of
 386 these biases provides an opportunity to address them more thoughtfully in biodiversity models
 387 that depend on these data. Specifically, we identified five collection norms common to the

388 practices of all collectors: they tend to collect a.) more species rather than multiple specimens of
389 the same species; b.) about 10 specimens per locality during their lifetime; c.) from localities
390 sampled by other collectors; d.) during the peak growing season in spring and summer when
391 climates are more favorable and photosynthetic rates and reproduction are generally higher; e.)
392 species from smaller genera and families; and f.) particular species that are available outside of
393 peak collecting months (i.e., when climates are less favorable for plant growth. We also
394 identified that some collections norms have changed through time with collectors avoiding
395 several taxonomically complex taxa during some decades.

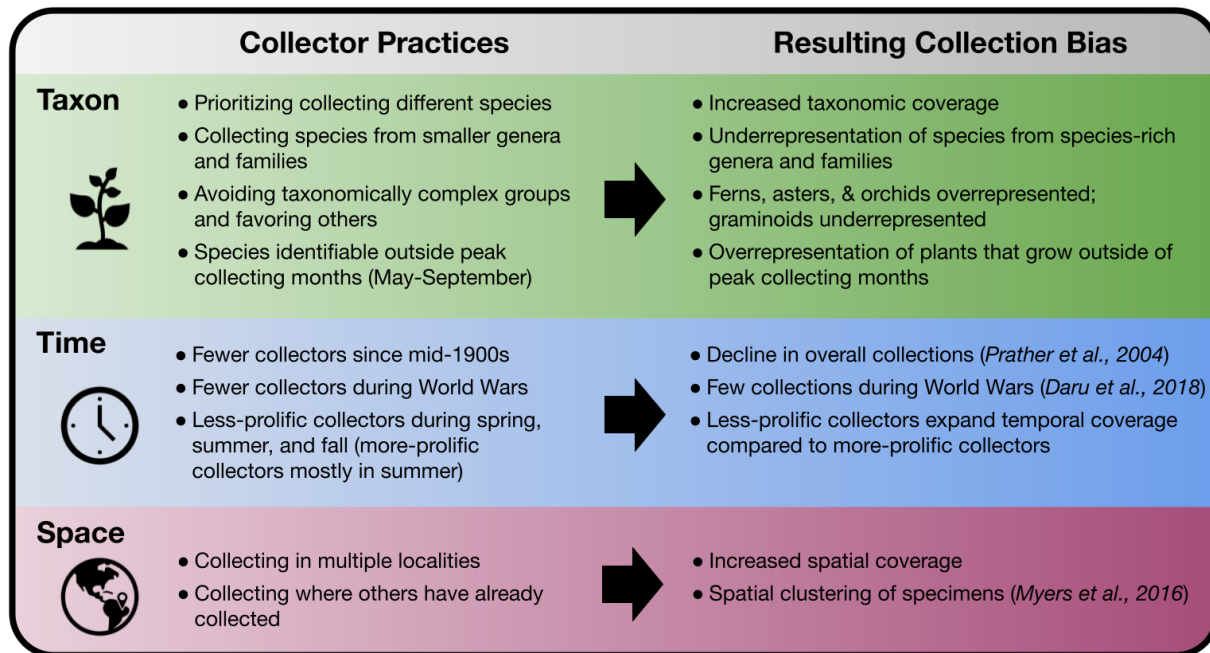
396

397 In contrast to the collections norms detailed above, we also identified several divergences
398 between the collector practices of more- versus less-prolific collectors. Specifically, more-prolific
399 collectors i.) collected largely during fewer months; ii.) had stronger affinities to certain localities;
400 and iii.) were not active in several large regions sampled by less-prolific collectors (e.g., the
401 state of New York, USA).

402

403 A summary of our findings is presented in Fig. 8, where we outline the collector practices and
404 resulting collection biases we have identified in the context of three key dimensions of bias:
405 taxon, time, and space. We include two previously identified temporal collection biases, the
406 decline in overall collections that was first presented by Prather *et al.* (2004) and the decline in
407 collections during World Wars I and II identified by Daru *et al.* (2018). We also include the
408 overall spatial clustering of collections, which was first defined by Myers *et al.* (2016).

409



410
 411 Figure 8. This graphic describes the collector practices that have shaped overall collection bias
 412 in natural history collections along three dimensions: taxon, time, and space.

413
 414 **Discussion**

415
 416 ***Taxonomic bias: prioritizing greater species diversity and the underrepresentation of***
 417 ***large, complex taxa***

418 We found that botanists in the northeastern US prioritized collecting more species versus
 419 collecting multiple specimens of the same species. Although this tendency has been viewed as
 420 problematic in biology (Lewin, 1982; May, 2004), we assert that such collecting has contributed
 421 considerably to expanding taxonomic coverage represented in NHCs. Despite this tendency,
 422 however, collectors do not sample species randomly: many collect the same taxa while avoiding
 423 others (Fig. 6). For instance, the brightly colored pink lady's slipper orchid (*Cypripedium acaule*
 424 Aiton) was collected by many people whereas hairy panicgrasses (*Dichanthelium acuminatum*
 425 (Sw.) Gould & C.A.Clark) was collected by relatively few. This collection norm affects our
 426 attempts to model biodiversity owing to the gap between taxon diversity and abundance
 427 information recorded in NHCs versus their actual diversities and abundances in nature (Elith &
 428 Leathwick, 2007; Gomes *et al.*, 2018). This pattern mirrors the collection norm whereby
 429 collectors tend to collect ten specimens per locality and suggests that collectors travelled to

430 different localities to collect new species rather than comprehensively collecting at a single
431 locality.
432
433 Taxonomic collection norms have likely contributed to the overrepresentation of less species-
434 rich taxa with distinctive morphologies (e.g., *Lobelia*, *Polystichum*, and *Dryopteris*) in herbaria
435 relative to larger taxa that are often taxonomically challenging (e.g., *Carex*, *Crataegus*, and
436 *Salix*). Specimens from many large taxa were collected by fewer people than expected,
437 suggesting these were mainly collected by botanists with specialized taxonomic interests. In the
438 northeastern US, such specialist-prone taxa include genera like *Sphagnum* (peat mosses),
439 *Dichanthelium* (rosette grasses), *Salix* (willows), and *Crataegus* (hawthorns), and families like
440 Poaceae, Cyperaceae, and Juncaceae (collectively, the graminoids). These groups often
441 require microscopic examination to distinguish subtle differences necessary for accurate
442 species identification and often can only be identified with reproductive features at specific
443 maturation stages (FNA Editorial Committee, eds., 1993+). Further complicating species
444 identification and delimitation are their complex evolutionary histories, including infrageneric
445 hybridization (Ennos *et al.*, 2005). We hypothesize that this taxonomic bias in collections is often
446 driven by the perceived taxonomic complexity and difficulty to identify species within such
447 groups (for discussions of taxonomic complexity, see Ennos *et al.*, 2005; Karbstein *et al.*, 2024).
448 This collection norm suggests that the most diverse groups, which are likely in greatest need of
449 study, are woefully underrepresented in NHCs.

450
451 We also identified clear trends in shifting taxonomic collection norms through time, a pattern that
452 has received little attention. We observed that taxonomic biases have apparently shifted, with
453 certain taxa being favored and others apparently avoided across different generations of
454 botanists. For example, in the northeastern US, many collectors in the 1930s avoided families
455 like Poaceae, Cyperaceae, Juncaceae, and Sphagnaceae. We hypothesize that collectors from
456 the Citizens Conservation Corps, many of whom lacked formal botanical training, may have
457 avoided families they perceived as more complex. In other words, we hypothesize that
458 collectors are less prone to collect what they don't know. This has significant implications for
459 comparing temporal trends between taxa; variations in historical collection intensity may affect
460 apparent changes in characteristics such as species distribution modeling (Franklin & Miller,
461 2009) and phenology (Miller-Rushing *et al.*, 2008). Therefore, understanding the overall
462 temporal distribution of collections is crucial for appreciating how record availability—and the
463 uncertainty in these data—changes over time.

464

465 ***Spatial bias: less-prolific collectors contribute unique spatial coverage with more-***
466 ***random spatial sampling***

467 We identified an important divergent collection practice between more- and less-prolific
468 collectors whereby less-prolific collectors contribute unique spatial coverage versus collections
469 by more-prolific collectors (see Fig. 4). These less-prolific collectors enhance sampling near
470 commonly collected localities and act as the backbone for entire regions where more-prolific
471 collectors have not collected, such as most of New York State, excluding New York City and
472 Long Island (see Fig. 5). Thus, the cumulative spatial coverage by more-prolific collectors is
473 considerably lower than that of less-prolific collectors, indicating that the collections made by the
474 latter more accurately reflect plant diversity across different regions.

475

476 Interestingly, the spatial bias of less-prolific collectors does not differ significantly from the
477 overall spatial bias in herbaria. However, these collections are still biased with respect to
478 random sampling. This suggests that while less-prolific collectors do not exhibit the same
479 preference for specific collection sites as more-prolific collectors, they also tend to revisit
480 locations where collections have previously been made. Despite this spatial collection norm, the
481 increased spatial coverage provided by less-prolific collectors has greatly improved the overall
482 spatial sampling in herbaria. This increased spatial coverage has helped facilitate the recent
483 application of herbarium data to disciplines that rely on extensive sampling; for example,
484 ecology (Meineke *et al.*, 2019a; Heberling, 2022); invasion biology (Crawford & Hoagland 2009;
485 Schmidt *et al.*, 2023), species distribution modeling (Daru *et al.*, 2021), environmental science
486 (Carbone *et al.*, 2023; Jakovljević *et al.*, 2024), and conservation biology (Schatz, 2002).

487

488 Finally, the broad spatial sampling by numerous less-prolific collectors that we identified reflects
489 patterns also observed with contemporary iNaturalist data, where contributions by millions of
490 community scientists greatly help extend spatial sampling beyond what is captured in herbaria
491 (Eckert *et al.*, 2024). This similarity indicates that the spatial biases of community scientists align
492 more closely with those of less-prolific collectors than with the more-prolific collectors who
493 contributed heavily to overall spatial biases in collections.

494

495 ***Temporal bias: variability driven by collector activity***

496 The substantial declines in collections over the past 75 years is consistent with trends observed
497 in other regional floras (Prather *et al.*, 2004; Daru *et al.*, 2018) and is strongly correlated with

498 declines in the number of active collectors. This suggests that while more-prolific collectors may
499 heavily influence the interannual intensity of collections at certain times (Bebber *et al.*, 2012;
500 Daru *et al.*, 2018), the overall trends are primarily driven by fluctuations in the number of all
501 active collectors.

502
503 Notably, the reduction in annual collections coincided with the years when the US was involved
504 in World Wars I (1917–1920) and II (1941–1946) when citizens from the northeastern US were
505 conscripted for military service. Following decreased collections during World War I, the spike in
506 collections and active collectors from 1932 through 1941 corresponds with US government
507 efforts to reduce unemployment and support environmental projects during the Great
508 Depression (1929–1939; Salmond, 1967). During this period, the government employed
509 thousands of citizens—primarily young men aged 18 to 25—for projects focusing on
510 environmental improvements (e.g., in the Civilian Conservation Corps; Salmond, 1967). A key
511 objective of these initiatives was to produce local species inventories, documented through
512 "complete herbaria," to aid in land planning and protection (Department of the Interior, 1936).
513 Since these projects often targeted similar habitats—primarily forested areas—many inventories
514 likely covered areas with similar species composition in the northeastern US. Consequently,
515 despite the spikes in collections, active collectors, and collection locations during this time, the
516 number of species collected during this period did not increase substantially. Once World War II
517 began and people from the same demographic were heavily drafted into WWII, all metrics once
518 again quickly declined. This highlights how major socio-political events affecting significant
519 population segments can directly impact NHCs by reducing the pool of available collectors.
520 Similar impacts of socio-political events on NHCs were recently documented in collection
521 requests for multiomic sampling, which plummeted during the global COVID pandemic (Davis *et*
522 *al.*, 2024).

523
524 We identified that less-prolific collectors increased overall sampling at the start and end of the
525 primary growing season (late spring and early autumn), which diverges from collections by
526 more-prolific collectors whose activity during these periods markedly decreases. The intensity of
527 sampling during these off-peak periods is crucial for improving the accuracy of phenological
528 estimates (Miller-Rushing *et al.*, 2008) and understanding the impact of anthropogenic climate
529 change on early- and late-season species (Kudo & Ida, 2013; Park *et al.*, 2023). We
530 hypothesize that the increased sampling by less-prolific collectors at the beginning and end of

531 the growing season (i.e., April–May and September–October) might be related to student
532 collections in university botany classes during the academic year (typically September–May).

533

534 Surprisingly, although 90% of specimens are collected in the northeastern US between May and
535 September, species collected outside the peak months are disproportionately represented
536 among the most abundant species in herbaria. These include many evergreen (e.g.,
537 *Polystichum acrostichoides* (Michx.) Schott and *Dryopteris marginalis* (L.) A.Gray), woody (e.g.,
538 *Vaccinium corymbosum* L. and *Acer rubrum* L.), and early-flowering species (e.g., *Viola sororia*
539 Willd. and *Arisaema triphyllum* (L.) Schott), as well as species with winter-available flowers or
540 fruits (e.g., *Ilex verticillata* (L.) A.Gray and *Hamamelis virginiana* L.). We hypothesize this
541 overrepresentation is driven by collectors' familiarity with these species, which are more
542 accessible and—in some cases—more identifiable outside of peak collection months when
543 fewer species are available.

544

545 ***Exceptions to the norms: unique collector practices of collectors contribute overall bias***

546 Despite the similar collector practices we identified, we emphasize that understanding how
547 some collectors diverged from these norms is important for understanding overall collection bias
548 in NHCs. For example, the most prolific collector in our dataset, R. L. Schaeffer, Jr., collected
549 50,287 specimens from only 195 localities—far fewer than expected based on our model. He
550 collected, almost exclusively, in the vicinity of Allentown, PA where Schaeffer taught botany at
551 Muhlenberg College from 1954-1983 ('R. L. Schaeffer Obituary', 2001). His singular efforts had
552 an outsized impact on overall spatial bias in the northeastern US with his collections being the
553 main driver of the high collection density in eastern PA, one of the most collection-dense areas
554 in the northeastern US. Furthermore, the expansive taxonomic coverage and high collection
555 depth of Schaeffer's specimens provides a rich documentation of the flora of eastern
556 Pennsylvania over nearly a half century that can be leveraged for a diversity of collections-
557 based investigations (e.g., Meineke *et al.*, 2019b). This highlights how integrating historical
558 information about collectors (especially mega-collectors like Schaeffer) can help explain the
559 more stochastic processes in biodiversity data and can illuminate important datasets better
560 characterizing species and ecosystem responses to anthropogenic pressures.

561

562 ***Conclusion***

563 Our findings reveal how our understanding of biodiversity is founded on the cumulative effort of
564 thousands of people, many of whom have made small but impactful contributions to natural

565 history collections (NHCs). The cumulative spatial, temporal, and taxonomic practices of all
566 collectors give rise to the overall biases in collections. It is crucial that we identify and categorize
567 these collector practices to better understand the drivers of overall collection bias in NHCs and
568 begin developing tools to address them. We have identified numerous predictable collection
569 norms that appear to have shaped overall bias in NHCs. The predictability of these biases
570 provides an exciting and promising opportunity to begin incorporating statistical tools to address
571 collection biases in biodiversity models. These results can also be leveraged to guide future
572 collection efforts that can minimize gaps in collections and reduce bias in NHCs moving forward.
573 We highlight that collector practices—even by those who collected only a small number of
574 specimens—have vastly expanded the coverage of NHCs and we assert that continued
575 collections of all sizes are crucial for continuing to expand the coverage of NHCs and further
576 increasing their utility for understanding biodiversity in the face of global change.

577

578 **Acknowledgements**

579 The authors thank Aaron Ellison for guidance on statistical analysis and comments on an earlier
580 draft of this manuscript. We also thank Jonnathan Kennedy for providing updates to the Harvard
581 Index of Botanists, Nawal Shrestha for support with coding and analysis, and the rest of the
582 Davis lab for input throughout the project. RJS and KES acknowledge support from the National
583 Science Foundation (Graduate Research Fellowship Program) and the Harvard University
584 Herbaria. CCD acknowledges funding from Harvard University and from National Science
585 Foundation funding grants: DEB 1754584, EF1208835, DEB 2101884, DEB 1802209, and MRA
586 2105903. LS acknowledges support from USDA Hatch project 1026539, NSF-DEB 1601101,
587 and School of Environmental and Biological Sciences to Chrysler Herbarium at Rutgers
588 University.

589

590 **Author Contribution**

591 RJS, CCD, and LS conceptualized the study. RJS and CCD developed the methodology, RJS
592 and KES led the data curation, and the investigations and formal analysis were completed by
593 RJS. RJS led data visualization with support from CCD, LS, and KES. RJS and CCD led writing
594 with input and support from LS and KES.

595

596 **Data Availability Statement**

597 The data generated during this study are available in the supporting information of this
598 manuscript. Table S2 (all georeferenced records used in this study) and all code created for this

599 study are available on Github ([DOI TO BE ADDED AFTER REVIEW, code available to
600 reviewers as an Rmd file]).

601

602 **Conflict of Interest Statement**

603 CCD declares that he is supported by LVMH Research and Dior Science, a company involved
604 in the research and development of cosmetic products based on floral extracts. He also serves
605 as a member of Dior's Age Reverse Board.

606

607 **References**

608 **Bebber DP, Carine MA, Davidse G, Harris DJ, Haston EM, Penn MG, Cafferty S, Wood JRI,**
609 **Scotland RW. 2012.** Big hitting collectors make massive and disproportionate contribution to
610 the discovery of plant species. *Proceedings of the Royal Society B: Biological Sciences* **279**:
611 2269–2274.

612 **Carbone MS, Ayers TJ, Ebert CH, Munson SM, Schuur EAG, Richardson AD. 2023.**
613 Atmospheric Radiocarbon for the Period 1910-2021 Recorded by Annual Plants. *Radiocarbon*
614 **65**: 357–374.

615 **CBD. 2022.** *Decision adopted by the conference of the parties to the convention on biological*
616 *diversity 15/4. Kunming-montreal global biodiversity framework.* Montreal, Canada.

617 **Crawford PHC, Hoagland BW. 2009.** Can herbarium records be used to map alien species
618 invasion and native species expansion over the past 100 years? *Journal of Biogeography* **36**:
619 651–661.

620 **CT DEEP. 2023.** Northeastern States Town Boundary Set. *Connecticut Department of Energy &*
621 *Environmental Protection.*

622 **Daru BH, Davies TJ, Willis CG, Meineke EK, Ronk A, Zobel M, Pärtel M, Antonelli A, Davis**
623 **CC. 2021.** Widespread homogenization of plant communities in the Anthropocene. *Nature*
624 *Communications* **12**: 6983.

625 **Daru BH, Park DS, Primack RB, Willis CG, Barrington DS, Whitfeld TJS, Seidler TG,**
626 **Sweeney PW, Foster DR, Ellison AM, et al. 2018.** Widespread sampling biases in herbaria
627 revealed from large-scale digitization. *New Phytologist* **217**: 939–955.

628 **Davis CC. 2023.** The herbarium of the future. *Trends in Ecology & Evolution* **38**: 412–423.

629 **Davis CC. 2024.** Collections are truly priceless. *Science* **383**: 1035–1035.

630 **Davis CC, Sessa E, Paton A, Antonelli A, Teisher JK. 2024.** Guidelines for the effective and
631 ethical sampling of herbaria. *Nature Ecology & Evolution*.

632 **Department of the Interior. 1936.** Annual Report of the Department of the Interior 1936.

633 **Elith J, Leathwick J. 2007.** Predicting species distributions from museum and herbarium
634 records using multiresponse models fitted with multivariate adaptive regression splines.
635 *Diversity and Distributions* **13**: 265–275.

636 **Ennos R, French G, Hollingsworth P. 2005.** Conserving taxonomic complexity. *Trends in*
637 *Ecology & Evolution* **20**: 164–168.

638 **Forest Inventory and Analysis. 2023.** U.S. Department of Agriculture, Forest Service,
639 *Research & Development*.

640 **Flora of North America Editorial Committee, eds. 1993.** Flora of North America (FNA). *Flora*
641 *of North America*.

642 **Franklin J, Miller JA. 2009.** *Mapping species distributions: spatial inference and prediction*.
643 Cambridge ; New York: Cambridge University Press.

644 **Funk VA. 2003.** 100 Uses for an Herbarium: well at least 72. *American Society of Plant*
645 *Taxonomists Newsletter* **17**: 17–19.

646 **GBIF.org. 2024.** GBIF Occurrence Download. Available at: <https://doi.org/10.15468/dl.rndw9f>.
647 (Accessed: 28 August 2024).

648 **Gomes VHF, IJff SD, Raes N, Amaral IL, Salomão RP, De Souza Coelho L, De Almeida**
649 **Matos FD, Castilho CV, De Andrade Lima Filho D, López DC, et al. 2018.** Species
650 Distribution Modelling: Contrasting presence-only models with plot abundance data. *Scientific*
651 *Reports* **8**: 1003.

652 **Groom Q, Bräuchler C, Cubey R, Dillen M, Huybrechts P, Kearney N, Klazenga N,**
653 **Leachman S, Paul DL, Rogers H, et al. 2022.** The disambiguation of people names in
654 biological collections. *Biodiversity Data Journal* **10**: e86089.

655 **Harvard University Herbaria. 2024.** Harvard Index of Botanists. *Harvard University Herbaria &*
656 *Libraries*.

657 **Heberling JM. 2022.** Herbaria as Big Data Sources of Plant Traits. *International Journal of*
658 *Plant Sciences* **183**: 87–118.

659 **Hedrick BP, Heberling JM, Meineke EK, Turner KG, Grassa CJ, Park DS, Kennedy J,**
660 **Clarke JA, Cook JA, Blackburn DC, et al. 2020.** Digitization and the Future of Natural History
661 Collections. *BioScience* **70**: 243–251.

662 **Jakovljević K, Mišljenović T, Van Der Ent A, Baker AJM, Invernón VR, Echevarria G. 2024.**
663 “Mining” the herbarium for hyperaccumulators: Discoveries of nickel and zinc
664 (hyper)accumulation in the genus *NOCCAEA* (Brassicaceae) through X-ray fluorescence
665 herbarium scanning. *Ecological Research* **39**: 450–459.

666 **Johnson KR, Owens IFP, the Global Collection Group. 2023.** A global approach for natural
667 history museum collections. *Science* **379**: 1192–1194.

668 **Karbstein K, Kösters L, Hodač L, Hofmann M, Hörandl E, Tomasello S, Wagner ND,**
669 **Emerson BC, Albach DC, Scheu S, et al. 2024.** Species delimitation 4.0: integrative taxonomy
670 meets artificial intelligence. *Trends in Ecology & Evolution* **39**: 771–784.

671 **Kozlov MV, Sokolova IV, Zverev V, Zvereva EL. 2021.** Changes in plant collection practices
672 from the 16th to 21st centuries: implications for the use of herbarium specimens in global
673 change research. *Annals of Botany* **127**: 865–873.

674 **Kudo G, Ida TY. 2013.** Early onset of spring increases the phenological mismatch between
675 plants and pollinators. *Ecology* **94**: 2311–2320.

676 **Lendemmer J, Thiers B, Monfils AK, Zaspel J, Ellwood ER, Bentley A, LeVan K, Bates J,**
677 **Jennings D, Contreras D, et al. 2020.** The Extended Specimen Network: A Strategy to
678 Enhance US Biodiversity Collections, Promote Research and Education. *BioScience* **70**: 23–30.

679 **Lewin R. 1982.** Biology Is Not Postage Stamp Collecting: Ernst Mayr, the eminent Harvard
680 evolutionist, explains why he thinks some physical scientists have a problem with evolution.
681 *Science* **216**: 718–720.

682 **Mancini M, Barber A, Block TA, Skema C. 2019.** Mid-Atlantic megalopolis georeferencing
683 guidelines.

684 **Marín-Rodulfo M, Rondinel-Mendoza KV, Martín-Girela I, Cañadas EM, Lorite J. 2024.** Old
685 meets new: Innovative and evolving uses of herbaria over time as revealed by a literature
686 review. *PLANTS, PEOPLE, PLANET* **6**: 1261–1271.

687 **May RM. 2004.** Tomorrow's taxonomy: collecting new species in the field will remain the rate–
688 limiting step (HCJ Godfray and S Knapp, Eds.). *Philosophical Transactions of the Royal Society*
689 *of London. Series B: Biological Sciences* **359**: 733–734.

690 **Meineke EK, Classen AT, Sanders NJ, Jonathan Davies T. 2019a.** Herbarium specimens
691 reveal increasing herbivory over the past century (A Iler, Ed.). *Journal of Ecology* **107**: 105–117.

692 **Meineke EK, Davies TJ, Daru BH, Davis CC. 2019b.** Biological collections for understanding
693 biodiversity in the Anthropocene. *Philosophical Transactions of the Royal Society B: Biological*
694 *Sciences* **374**: 20170386.

695 **Miller-Rushing AJ, Inouye DW, Primack RB. 2008.** How well do first flowering dates measure
696 plant responses to climate change? The effects of population size and sampling frequency.
697 *Journal of Ecology* **96**: 1289–1296.

698 **Park DS, Davis CC. 2017.** Implications and alternatives of assigning climate data to
699 geographical centroids. *Journal of Biogeography* **44**: 2188–2198.

700 **Park DS, Xie Y, Ellison AM, Lyra GM, Davis CC. 2023.** Complex climate-mediated effects of
701 urbanization on plant reproductive phenology and frost risk. *New Phytologist* **239**: 2153–2165.

702 **Pebesma E. 2018.** Simple Features for R: Standardized Support for Spatial Vector Data. *The R*
703 *Journal* **10**: 439.

704 **Pebesma E, Bivand R. 2023.** *Spatial Data Science: With Applications in R*. New York:
705 Chapman and Hall/CRC.

706 **PennDOT. 2024.** PennDOT Open Data. *Pennsylvania Department of Transportation*.

707 **Perring FH, Walters SM (Eds.). 1962.** *Atlas of the British flora*. London: Thomas Nelson &
708 Sons.

709 **Prather LA, Alvarez-Fuentes O, Mayfield MH, Ferguson CJ. 2004.** The Decline of Plant
710 Collecting in the United States: A Threat to the Infrastructure of Biodiversity Studies. *Systematic*
711 *Botany* **29**: 15–28.

712 **Preston CD. 2013.** Following the BSBI's lead: the influence of the *Atlas of the British flora* ,
713 1962–2012. *New Journal of Botany* **3**: 2–14.

714 **Obituary for Robert L. Schaeffer (Aged 83). 2001.** *The Morning Call*: 28.

715 **Rudis VA. 2003.** *Comprehensive Regional Resource Assessments and Multipurpose Uses of*
716 *Forest Inventory and Analysis Data, 1976 to 2001: A Review.* Asheville, North Carolina: USDA
717 Forest Service, Southern Research Station.

718 **Salmond JA. 1967.** *The Civilian Conservation Corps, 1933-1942; a New Deal case study.*
719 Durham, North Carolina: Duke University Press.

720 **Schatz GE. 2002.** Taxonomy and Herbaria in Service of Plant Conservation: Lessons from
721 Madagascar's Endemic Families. *Annals of the Missouri Botanical Garden* **89**: 145.

722 **Schmidt RJ, King MR, Aronson MFJ, Struwe L. 2023.** Hidden cargo: The impact of historical
723 shipping trade on the recent-past and contemporary non-native flora of northeastern United
724 States. *American Journal of Botany* **110**: e16224.

725 **Schorn C, Weber E, Bernardos R, Hopkins C, Davis C. 2016.** The New England Vascular
726 Plants Project: 295,000 specimens and counting. *Rhodora* **118**: 324–325.

727 **Shorthouse DP. 2024.** Bionomia. *Bionomia*.

728 **Silge J, Robinson D. 2016.** tidytext: Text Mining and Analysis Using Tidy Data Principles in R.
729 *The Journal of Open Source Software* **1**: 37.

730 **Sweeney PW, Starly B, Morris PJ, Xu Y, Jones A, Radhakrishnan S, Grassa CJ, Davis CC.**
731 **2018.** Large-scale digitization of herbarium specimens: Development and usage of an
732 automated, high-throughput conveyor system. *TAXON* **67**: 165–178.

733 **United States Census Bureau. 2024.** Cartographic Boundary Files: States: 1 : 500,000
734 (national). Available at: [https://www.census.gov/geographies/mapping-files/time-](https://www.census.gov/geographies/mapping-files/time-series/geo/cartographic-boundary.html)
735 [series/geo/cartographic-boundary.html](https://www.census.gov/geographies/mapping-files/time-series/geo/cartographic-boundary.html). (Accessed: 6 September 2024).
736

737 **Webster MS (Ed.). 2017.** *The extended specimen: emerging frontiers in collections-based*
738 *Ornithological Research.* Boca Raton London New York: CRC Press, Taylor & Francis Group.

739 **Weeks A, Collins E, Majors T, Murrell Z, Paul D, Sheik M, Shorthouse D, Zeringue-**
740 **Krosnick S. 2024.** Workshop Report: Supporting inclusive and sustainable collections-based
741 research infrastructure for systematics (SISRIS). *Research Ideas and Outcomes* **10**: e126532.

742 **Whelan A. 1948.** Of People and Places. *The Bridgeport Sunday Post*: 23.

743

744 **Supporting Information**

745 **Table S1** Herbaria whose specimens were used for this study, indicating the institution code,
746 institution name, and the number of specimens from each herbarium that were used in this
747 study.

748 **Table S2** Total specimens used in this study after data cleaning, georeferencing, and collector
749 disambiguation.

750 **Table S3** A table containing the DarwinCore recordedBy strings from gbif and the unique
751 identifier representing each collector.

752 **Table S4** The one hundred most frequently collected species in the northeastern US.

753 **Fig. S1** A boxplot showing the difference in number of specimens of each species related to
754 whether the species has been collected only during peak collection months (May, June, July,
755 August, and September) or also collected in non-peak months.

***New Phytologist* Supporting Information**

Article title: Identifying the collector practices that shape spatial, temporal, and taxonomic bias in herbaria

Authors: Ryan J. Schmidt, Kristen E. Saban, Lena Struwe, Charles C. Davis

Article acceptance date: **TBD**

The following Supporting Information is available for this article:

Table S1 Herbaria whose specimens were used for this study, indicating the institution code, institution name, and the number of specimens from each that were used in this study.

Table S2 Total specimens used in this study after data cleaning, georeferencing, and collector disambiguation.

Table S3 A table containing the DarwinCore recordedBy strings from gbif and the unique identifier representing each collector.

Table S4 The 100 most frequently collected species in the northeastern US.

Fig. S1 A boxplot showing difference in number of specimens of each species related to whether the species has been collected only during peak collection months (May, June, July, August, and September) or also collected in non-peak months.

institutionCode	organization	specimens
PH	Academy of Natural Sciences	167758
NEBC	New England Botanical Club	167252
CM	Carnegie Museum of Natural History	154796
NY	The New York Botanical Garden	102564
UCONN	George Safford Torrey Herbarium, University of Connecticut	88196
YPM (YU)**	Yale Peabody Museum	51292
MCA	Muhlenberg College	40003
GH	Harvard University	38760
VT	University of Vermont, Plant Biology	19399
MOAR	Morris Arboretum, University of Pennsylvania	17888
BUF	Buffalo Museum of Science	13081
F	Field Museum of Natural History	10103
MVSC	Millersville University	8204
SIM	Staten Island Museum	7829
A	Harvard University	7036
CHRB	Rutgers University	6249
MICH	University of Michigan	5515
DUKE	Duke University	4995
USF	University of South Florida	4933
KU	Kwangsi University	4029
OSW**	State University of New York at Oswego	3710
EIU	Eastern Illinois University	2623
NCU	University of North Carolina at Chapel Hill	2377
WVW	West Virginia Wesleyan College	2338
TENN	University of Tennessee - Knoxville	2262
CAS	California Academy of Sciences	2037
BDI	Putnam Museum of History and Natural Science	1935
BRIT	Botanical Research Institute of Texas	1854
ECON	Harvard University	1638
DOV	Delaware State University	1603
PRC	Charles University, Prague	1585
CMMF**	Université de Montréal Biodiversity Centre	1401

BRY	Brigham Young University	1356
FH	Harvard University	1327
NCSC	North Carolina State University	1254
US	Smithsonian Institution	1248
AMES	Harvard University	1232
CGCC	Columbia-Greene Community College	1207
MIN	University of Minnesota	1193
LSU	Louisiana State University	1188
UTEP	University of Texas at El Paso	1147
UMO	University of Missouri	1073
MU	Miami University	969
SDSU	San Diego State University	966
MPM**	Milwaukee Public Museum	954
RSA	California Botanic Garden	949
FLAS	Florida Museum of Natural History	825
USU**	Utah State University	812
MISS	University of Mississippi	785
CHAS	Southern Research Station, USDA Forest Service	762
DEK	Northern Illinois University	707
UCR	University of California, Riverside	700
IAC	Instituto Agronômico de Campinas	699
WS	Washington State University	687
MISSA	Mississippi State University	682
SD	San Diego Natural History Museum	648
SBBG	Santa Barbara Botanic Garden	602
BBM**	Beaty Biodiversity Museum, University of British Columbia	598
ASU	Arizona State University	569
IBUNAM*	National Autonomous University of Mexico Herbarium	551
COLO	University of Colorado Museum of Natural History	535
DBG	Denver Botanic Gardens	519
NHA	University of New Hampshire	519
AUA**	John D. Freeman Herbarium, Auburn University Museum of Natural History	506

TAES	Texas A&M University	504
ALTA/UADBG**	University of Alberta Museums	501
MSC	Michigan State University	487
APCR	Arkansas Tech University	456
MWI	Western Illinois University	436
CHSC	California State University, Chico	413
NO	Tulane University	401
HUDC	Howard University	370
WVA	West Virginia University	369
KSP	Pittsburg State University	368
CINC	University of Cincinnati	356
SAT	Angelo State University	352
MO	Missouri Botanical Garden	348
SFV	California State University, Northridge	330
ISC	Iowa State University	315
LD	Lund University	314
CLEMS	Clemson University	310
LOB	California State University, Long Beach Herbarium	288
OS	Ohio State University	278
MEL	Royal Botanic Gardens Victoria	264
IDS	Idaho State University	263
FTG	Fairchild Tropical Botanic Garden	261
UNM	University of New Mexico	242
ROM (TRT/TRTC)**	Royal Ontario Museum	240
UdeM**	Université de Montréal	237
UT	University of Utah	223
MUHW	Marshall University	222
FSU	Florida State University	212
CDA	California Department of Food and Agriculture	207
GREE	University of Northern Colorado	181
CS	Colorado State University	179
GA	University of Georgia	179
UWW	University of Wisconsin - Whitewater	174

MMNS	Mississippi Museum of Natural Science	172
DSRC*	Mohonk Preserve	157
JSNM	Jurica-Suchy Nature Museum at Benedictine University	155
NEON	Arizona State University	149
LA	University of California, Los Angeles	138
MA	Real Jardín Botánico	131
UNA	University of Alabama	131
ID	University of Idaho	129
UWMB(WTU)**	University of Washington	126
UMD (MARY)**	University of Maryland	123
BAYLU	Baylor University	118
SRP	Boise State University	114
USCH	University of South Carolina	113
TRTE*	University of Toronto Mississauga	112
NCSM	North Carolina Museum of Natural Sciences	112
UFPR (UPCB)**	Universidade Federal do Paraná	108
ENCB-IPN (ENCB)**	Instituto Politécnico Nacional	102
NHMUK (BM)**	The Natural History Museum	101
OSU (OSUF)**	Oregon State University	90
HTTU	Tennessee Technological University	85
OBI	California Polytechnic State University	81
FSC	California State University, Fresno	78
TTC	Texas Tech University	77
EWU	Eastern Washington University	73
DES	Desert Botanical Garden	70
POM	Pomona College	66
NBM	New Brunswick Museum	64
UAM	University of Arkansas at Monticello	62
MSUB	Montana State University-Billings	60
RENO	University of Nevada	58
TRH	Norwegian University of Science and Technology	56
NMNZ*	New Zealand National Museum of Natural History	53
KUN	Kunming Institute of Botany, Chinese Academy of Sciences	49

SJSU	San Jose State University	49
CAU	Campbell University	48
BMO	Unknown	47
ODU	Old Dominion University	45
HPSU	Portland State University	45
EKY	Eastern Kentucky University	42
BOON	Appalachian State University	41
IUP	Indiana University of Pennsylvania	41
ACAD	Acadia University	37
PUA	Pacific Union College	32
HO	Tasmanian Museum and Art Gallery	32
WCW	Whitman College	32
UARK	University of Arkansas	31
NEB	University of Nebraska State Museum, Lincoln NE	31
SMU	Southern Methodist University	30
WCUH	Western Carolina University	30
WSCO	Weber State University	28
UVSC	Utah Valley University	25
SAU	Sichuan Agricultural University	23
UCSB	University of California, Santa Barbara	23
BUT	Butler University	21
MACF	California State University Fullerton	19
JBRJ (RB)**	Rio de Janeiro Botanical Garden herbarium	19
ASC	Northern Arizona University	17
WWB	Western Washington University	17
UWL	University of Wisconsin	16
CUP	Cornell University	15
GMUF	George Mason University	15
H	University of Helsinki	15
CSLA	California State University	14
EMC	Eastern Michigan University	14
MPEG (MG)**	Museu Paraense Emílio Goeldi	14
BC	Institut Botànic de Barcelona	13

LFCC	Lord Fairfax Community College	13
TAWES	Maryland Department of Natural Resources	12
IRVC	University of California, Irvine	12
CIC	The College of Idaho	11
BEREA	Berea College	10
PSM	Slater Museum of Natural History, University of Puget Sound	10
USMS	University of Southern Mississippi	10
DAV	University of California, Davis	9
CONN	University of Connecticut	9
NYS	New York State Museum	8
SOC	Southern Oregon University	8
BRU	Brown University	7
COCC	Central Oregon Community College	7
UNESP-FCA	Unkown	7
UAC	University of Calgary	7
LEA	University of Lethbridge	7
AU	Xiamen University	7
BABY	Yukon Government	7
BH	Cornell University	6
HSC	Humboldt State University	6
SHM	Shanghai Museum of Natural History	6
KSTC	Emporia State University	5
IBE	Institute for Botanical Exploration	5
GINCO	Agriculture & Agri-Food Canada	5
RBGE/E**	Royal Botanical Gardens Edinburgh	5
Royal Botanical Gardens	Unknown	5
Utah Tech University*	Utah Tech University	5
NFLD/SWGC**	Memorial University of Newfoundland	4
UCSC	University of California Santa Cruz	4
ETSU	East Tennessee State University	3
TAM	Estonian Museum of Natural History	3

IND	Indiana University	3
MBM	Museu Botânico Municipal	3
BRFC*	Black Rock Forest Consortium Herbarium	3
CNS-UT (CNS)**	Australian Tropical Herbarium	3
MOR	The Morton Arboretum	3
MASS	University of Massachusetts	3
UWO	University of Western Ontario	3
SNM	Western New Mexico University	3
CSUSB	California State University, San Bernardino	2
NAS	Institute of Botany, Jiangsu Province and Chinese Academy of Sciences	2
AAFC	National Collection of Vascular Plants, Agriculture and Agri-Food Canada	2
Unknown	Unknown	2
OKLA	Oklahoma State University	2
VPI	Virginia Polytechnic Institute and State University	2
ANHC	Arkansas Natural Heritage Commission	1
CU	Cornell University	1
MACB	Facultad de Ciencias Biológicas, Universidad Complutense de Madrid	1
CORD	Herbario CORD	1
USZ	Herbario del Oriente Boliviano (USZ), Museo de Historia Natural Noel Kempff Mercado, UAGRM	1
PE	Institute of Botany, Chinese Academy of Sciences	1
LE	Komarov Botanical Institute of RAS	1
MNHN	Museo Nacional de Historia Natural	1
R	Museu Nacional	1
CIIDIR-IPN (CIIDIR)**	Instituto Politécnico Nacional, CIIDIR Unidad Durango	1
FML	Unknown	1
GenBank	GenBank	1
PSUC	Unknown	1
UACH	Unknown	1
UMKC*	University of Missouri - Kansas City	1

VALE	Unknown	1
W	Naturhistorisches Museum Wien	1
SFSU	San Francisco State University	1
FR	Senckenberg Gesellschaft für Naturforschung: Senckenberg Forschungsinstitut und Naturmuseum	1
M	Staatliche Naturwissenschaftliche Sammlungen Bayerns (SNSB)	1
BING	State University of New York	1
S	Swedish Museum of Natural History	1
TROM	UiT The Arctic University of Norway	1
USP	Universidad San Pablo-CEU	1
UESC	Universidade Estadual de Santa Cruz	1
MONTU	University of Montana	1
UNB	University of New Brunswick	1
OULU	University of Oulu	1
TEX	University of Texas at Austin	1

Table S1 Herbaria whose specimens were used for this study, indicating the institutionCode from gbif, the institution name from Index Herbariorum (<https://sweetgum.nybg.org/science/ih/>), and the number of specimens from each herbarium that were used in this study.

* Herbaria that are not included in Index Herbariorum

** Herbaria that are listed under a different name in the gbif dataset and Index Herbariorum

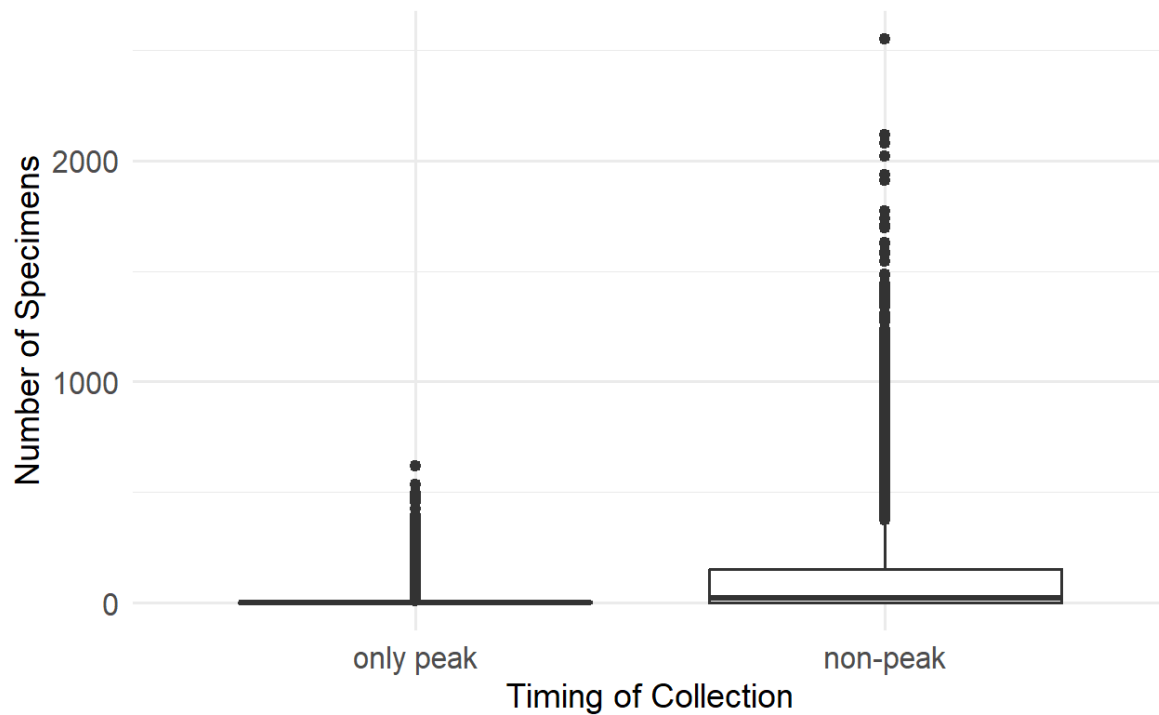


Fig. S1 A boxplot showing difference in number of specimens of each species related to whether the species has been collected only during peak collection months (May, June, July, August, and September) or also collected in non-peak months.

[Attached separately]

Table S2 Total specimens used in this study after data cleaning, georeferencing, and collector disambiguation.

[Attached separately]

Table S3 A table containing the DarwinCore recordedBy strings from gbif and the unique identifier representing each collector.

#	family	acceptedScientificName	specimens
1	Ericaceae	<i>Vaccinium corymbosum</i> L.	2597
2	Dryopteridaceae	<i>Polystichum acrostichoides</i> (Michx.) Schott	2231
3	Dryopteridaceae	<i>Dryopteris intermedia</i> (Muhl. ex Willd.) A.Gray	2139
4	Violaceae	<i>Viola sororia</i> Willd.	2056
5	Ophioglossaceae	<i>Sceptridium dissectum</i> (Spreng.) Lyon	2023
6	Dryopteridaceae	<i>Dryopteris marginalis</i> (L.) A.Gray	2022
7	Athyriaceae	<i>Athyrium angustum</i> (Willd.) C.Presl	1863
8	Araceae	<i>Arisaema triphyllum</i> (L.) Schott	1774
9	Aquifoliaceae	<i>Ilex verticillata</i> (L.) A.Gray	1732
10	Cyperaceae	<i>Scirpus cyperinus</i> (L.) Kunth	1715
11	Dryopteridaceae	<i>Dryopteris carthusiana</i> (Vill.) H.P.Fuchs	1694
12	Equisetaceae	<i>Equisetum arvense</i> L.	1670
13	Dennstaedtiaceae	<i>Sitobolium punctilobum</i> (Poir.) Desv.	1659
14	Onocleaceae	<i>Onoclea sensibilis</i> L.	1654
15	Cyperaceae	<i>Carex lurida</i> Wahlenb.	1569
16	Thelypteridaceae	<i>Amauropelta noveboracensis</i> (L.) S.E.Fawc. & A.R.Sm.	1536
17	Cyperaceae	<i>Carex vulpinoidea</i> Michx.	1510
18	Ericaceae	<i>Gaylussacia baccata</i> (Wangenh.) K.Koch	1506
19	Osmundaceae	<i>Claytosmunda claytoniana</i> (L.) Metzgar & Rouhan	1501
20	Cyperaceae	<i>Carex scoparia</i> Schkuhr ex Willd.	1455
21	Sapindaceae	<i>Acer rubrum</i> L.	1454
22	Ophioglossaceae	<i>Botrypus virginianus</i> (L.) Michx.	1431
23	Viburnaceae	<i>Viburnum acerifolium</i> L.	1429
24	Hamamelidaceae	<i>Hamamelis virginiana</i> L.	1420
25	Oxalidaceae	<i>Oxalis stricta</i> L.	1411
26	Rosaceae	<i>Rubus allegheniensis</i> Porter	1389
27	Ericaceae	<i>Vaccinium pallidum</i> Aiton	1375
28	Cyperaceae	<i>Carex rosea</i> Willd.	1371
29	Poaceae	<i>Dichanthelium acuminatum</i> (Sw.) Gould & C.A.Clark	1344
30	Violaceae	<i>Viola cucullata</i> Aiton	1340
31	Ranunculaceae	<i>Thalictrum pubescens</i> Pursh	1339
32	Aspleniaceae	<i>Asplenium platyneuron</i> (L.) Britton, Sterns & Poggenb.	1337

33	Ericaceae	<i>Vaccinium angustifolium</i> Aiton	1331
34	Asteraceae	<i>Solidago juncea</i> Aiton	1320
35	Brassicaceae	<i>Cardamine pensylvanica</i> Muhl.	1312
36	Asteraceae	<i>Symphyotrichum lateriflorum</i> (L.) Á.Löve & D.Löve	1303
37	Asteraceae	<i>Solidago caesia</i> L.	1300
38	Cyperaceae	<i>Cyperus strigosus</i> L.	1290
39	Lycopodiaceae	<i>Diphasiastrum digitatum</i> (Dill. ex A.Braun) Holub	1286
40	Asteraceae	<i>Achillea millefolium</i> L.	1277
41	Pteridaceae	<i>Adiantum pedatum</i> L.	1268
42	Asteraceae	<i>Eurybia macrophylla</i> (L.) Cass.	1262
43	Oleaceae	<i>Fraxinus americana</i> L.	1262
44	Alismataceae	<i>Sagittaria latifolia</i> Willd.	1262
45	Lycopodiaceae	<i>Huperzia lucidula</i> (Michx.) Trevis.	1259
46	Juncaceae	<i>Juncus tenuis</i> Willd.	1259
47	Asteraceae	<i>Solidago nemoralis</i> Aiton	1257
48	Ranunculaceae	<i>Ranunculus abortivus</i> L.	1232
49	Dryopteridaceae	<i>Dryopteris cristata</i> (L.) A.Gray	1223
50	Polypodiaceae	<i>Polypodium virginianum</i> L.	1221
51	Asteraceae	<i>Symphyotrichum cordifolium</i> (L.) G.L.Nesom	1221
52	Euphorbiaceae	<i>Euphorbia maculata</i> L.	1214
53	Cyperaceae	<i>Carex laxiflora</i> Lam.	1213
54	Asteraceae	<i>Solidago rugosa</i> Mill.	1212
55	Balsaminaceae	<i>Impatiens capensis</i> Meerb.	1209
56	Violaceae	<i>Viola blanda</i> Willd.	1208
57	Ericaceae	<i>Gaultheria procumbens</i> L.	1201
58	Cyperaceae	<i>Carex intumescens</i> Rudge	1194
59	Cornaceae	<i>Cornus amomum</i> Mill.	1193
60	Asteraceae	<i>Solidago bicolor</i> L.	1180
61	Osmundaceae	<i>Osmundastrum cinnamomeum</i> subsp. <i>cinnamomeum</i>	1175
62	Asteraceae	<i>Eupatorium perfoliatum</i> L.	1168
63	Asteraceae	<i>Solidago gigantea</i> Aiton	1148
64	Lauraceae	<i>Lindera benzoin</i> (L.) Blume	1146
65	Asteraceae	<i>Antennaria plantaginifolia</i> (L.) Hook.	1145

66	Lycopodiaceae	<i>Dendrolycopodium obscurum</i> (L.) A.Haines	1134
67	Campanulaceae	<i>Lobelia inflata</i> L.	1129
68	Rubiaceae	<i>Mitchella repens</i> L.	1125
69	Rubiaceae	<i>Galium triflorum</i> Michx.	1118
70	Rosaceae	<i>Geum canadense</i> Jacq.	1117
71	Violaceae	<i>Viola pubescens</i> Aiton	1117
72	Viburnaceae	<i>Sambucus canadensis</i> L.	1116
73	Rosaceae	<i>Prunus serotina</i> Ehrh.	1115
74	Rosaceae	<i>Potentilla simplex</i> Michx.	1101
75	Cornaceae	<i>Cornus florida</i> L.	1098
76	Salicaceae	<i>Salix discolor</i> Muhl.	1098
77	Orchidaceae	<i>Cypripedium acaule</i> Aiton	1084
78	Cyperaceae	<i>Carex swanii</i> (Fernald) Mack.	1083
79	Salicaceae	<i>Salix eriocephala</i> Michx.	1079
80	Lauraceae	<i>Sassafras albidum</i> (Nutt.) Nees	1078
81	Araliaceae	<i>Aralia nudicaulis</i> L.	1075
82	Lamiaceae	<i>Lycopus americanus</i> Muhl. ex W.P.C.Barton	1065
83	Geraniaceae	<i>Geranium maculatum</i> L.	1059
84	Aristolochiaceae	<i>Asarum canadense</i> L.	1057
85	Sapindaceae	<i>Acer pensylvanicum</i> L.	1041
86	Ericaceae	<i>Kalmia angustifolia</i> L.	1039
87	Asteraceae	<i>Lactuca canadensis</i> L.	1036
88	Rosaceae	<i>Prunus virginiana</i> L.	1035
89	Ericaceae	<i>Rhododendron viscosum</i> (L.) Torr.	1034
90	Asparagaceae	<i>Maianthemum racemosum</i> (L.) Link	1033
91	Orchidaceae	<i>Spiranthes cernua</i> (L.) Rich.	1030
92	Lamiaceae	<i>Glechoma hederacea</i> L.	1026
93	Ericaceae	<i>Kalmia latifolia</i> L.	1026
94	Asteraceae	<i>Antennaria howellii</i> subsp. <i>neodioica</i> (Greene) R.J.Bayer	1023
95	Cyperaceae	<i>Eleocharis obtusa</i> (Willd.) Schult.	1022
96	Orchidaceae	<i>Goodyera pubescens</i> (Willd.) R.Br.	1020
97	Asparagaceae	<i>Maianthemum canadense</i> Desf.	1015

98	Rosaceae	<i>Rubus hispidus</i> L.	1011
99	Rosaceae	<i>Crataegus macrosperma</i> Ashe	1008
100	Lycopodiaceae	<i>Lycopodium clavatum</i> L.	1007

Table S4 The one hundred most frequently collected species in the northeastern US, including the family, scientific name (from gbif's acceptedScientificName field), and the number of specimens.