1    **The collector practices that shape spatial, temporal, and taxonomic bias in herbaria**

2

3    Ryan J. Schmidt[1]*, Kristen E. Saban[1], Lena Struwe[2,3], Charles C. Davis[1]*

4

5    [1] Department of Organismic and Evolutionary Biology, Harvard University Herbaria, Harvard
6        University, 22 Divinity Avenue, Cambridge, MA 02138, USA
7    [2] Department of Ecology, Evolution & Natural Resources, Rutgers, The State University of New
8        Jersey, 14 College Farm Road, New Brunswick, NJ 08901-8551, USA
9    [3] Department of Plant Biology, Rutgers, The State University of New Jersey, 59 Dudley Road,
10       New Brunswick, NJ 08901-8551, USA
11   * Corresponding authors: ryanschmidt@g.harvard.edu, cdavis@oeb.harvard.edu

12

13   RJS: 0000-0002-4907-2270
14   KES: 0009-0009-5728-4001
15   LS: 0000-0001-6074-5758
16   CCD: 0000-0001-8747-1101

17

18   **Word Count: 6334**

19   Introduction: 730
20   Materials & Methods: 1365
21   Results: 2115
22   Discussion: 2124

23

24   **Figures:**

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

27

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

29

## Summary

30

31 - Natural history collections (NHCs) are essential for studying biodiversity. Although
32   spatial, temporal, and taxonomic biases in NHCs affect analyses, the influence of
33   collector practices on biases remains largely unexplored.
34 - We utilized one million digitized specimens collected in the northeastern United States
35   by ~10,000 collectors to investigate (a) how collector practices shape spatial, temporal,
36   and taxonomic biases in NHCs and (b) similarities and differences between practices of
37   more- and less-prolific collectors
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

## Keywords

51

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

54

## Introduction

55

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 know 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 provide a representative sample of global biodiversity across
71    time, space, and taxa. Any deviations between a spatially, temporally, and taxonomically
72    representative sample and the representation of biodiversity in NHCs are examples of collection
73    bias. Understanding how NHCs diverge from this ideal coverage allows us to better account for
74    biases in our biodiversity models and discern what questions we can address using these
75    collections. Ultimately, understanding collection biases will help guide the application and
76    development of statistical tools to correct for biases, develop better priorities for future collecting
77    efforts, and help us achieve more comprehensive and accurate models of global biodiversity.
78
79    Comprehensive digitization of natural history specimens from large geographic/floristic regions
80    has revealed key spatial, temporal and taxonomic biases in NHCs (Meyer *et al.*, 2016; Daru *et*
81    *al.*, 2018; Kozlov *et al.*, 2021; Eckert *et al.*, 2024). These overall biases in NHCs are a
82    consequence of the spatial, temporal, and taxonomic collection practices of each collector—
83    what we call collector practices. Previous studies have highlighted the connection between
84    collector practices and overall bias in collections, documenting that a small number of mega-
85    collectors have made disproportionately large contributions to species discovery (Bebber *et al.*
86    2012) and to specimen collections in NHCs (Daru *et al.* 2018). The disproportionately large
87    impact of these mega-collectors raises an important but unanswered question: have highly
88    prolific collectors also contributed disproportionately to the biases documented in these
89    collections? To date, there have been no efforts to investigate how the collector practices of all
90    collectors in a region have contributed to overall bias in NHCs. Moreover, there have been no
91    large-scale efforts to understand the impact that less-prolific collectors have had on the spatial,
92    temporal, and taxonomic coverage in collections.
93
94    Here, we expand the current framework for investigating biases in NHCs (*sensu* Daru *et al.*,
95    2018) by explicitly examining how collection biases are shaped by the practices of individual
96    collectors which, to our knowledge, has not been broadly examined. . As a test case for our

97    investigation, we leverage the nearly completely digitized metaherbarium that extensively

98    documents the flora of the northeastern United States (i.e., all digitally available specimens

99    collected in the northeastern US and housed throughout the world; Schorn *et al.*, 2016;

100   Sweeney *et al.*, 2018; Hedrick *et al.*, 2020). Specifically, we use all digitized herbarium

101   specimens of land plants (i.e., bryophytes and vascular plants) collected in the northeastern

102   United States from the earliest digitized record to the present (i.e., 1781–2024). We reconstruct

103   the contributions of collectors to investigate how overall bias in NHCs are shaped by the

104   similarities and differences in collection practices of different collectors. We assess the

105   relationship between these collection practices and the number of collections by each collector

106   on a continuous scale with more- and less-prolific collectors representing opposite ends of this

107   continuum. Mega-collectors—who have contributed a disproportionately large amount of

108   specimens (*sensu* Daru *et al.*, 2018)—represent the uppermost extreme of this spectrum. We

109   also investigate how what we term *collection norms*—the collector practices shared by all

110   collectors—have influenced overall biases in NHCs. Such synthetic investigations further

111   demonstrate the growing utility of digitized specimens within the framework of the extended

112   specimen (Webster, 2017; Lendemer *et al.*, 2020), facilitating proper attribution for the

113   thousands of hidden heroes that have made meaningful but previously unrecognized

114   contributions to NHCs (Groom *et al.*, 2022) and enabling ongoing efforts to better model

115   biodiversity in an era of rapid ecological change.

116

117   **Materials and Methods**

118   ***Data collection & data cleaning***

119   We downloaded 2,365,287 records representing all digitized herbarium specimens of land

120   plants from the northeastern United States (i.e., Connecticut, Maine, Massachusetts, New

121   Hampshire, New Jersey, New York, Pennsylvania, Rhode Island, and Vermont; hereafter the

122   Northeast) from GBIF (GBIF.org, 2024). These specimens are housed in 237 herbaria around

123   the world (Table S1). We then filtered this dataset to remove the 548,895 records without a

124   transcribed date, collector, locality, or species-level identification. This filtering left us with

125   1,816,392 analyzable records.

126

127   ***Georeferencing***

128   About half of the cleaned records (920,633 records) contained transcribed coordinates. We

129   batch-georeferenced an additional 401,450 specimens to municipal centroid points (i.e., the

130   centroid points for local incorporated communities such as cities, towns, and townships; CT

131    DEEP, 2023; PennDOT, 2024) and removed all records that could not be georeferenced to a

132    specific municipality (503,563 records removed). Although this method of georeferencing does

133    not capture fine-scale differences in collection localities (Park & Davis, 2017), it is consistent

134    with the precision for many herbarium georeferencing initiatives in the northeastern US (e.g.,

135    Mancini *et al.*, 2019) and suitable for analyses on these large spatial scales. We removed

136    records with coordinates outside of the northeastern US (United States Census Bureau, 2024)

137    using the `st_intersection()` function from the *sf* package in R 4.4.1 (Pebesma 2018;

138    Pebesma & Bivand 2023; 10,254 records removed). This resulted in a total of 1,311,829

139    georeferenced records (see Fig. S1 for more information about the specimens removed at each

140    step of data cleaning).

141

142    ***Collector disambiguation***

143    Due to institutional differences in transcription practices, incorrect transcriptions, and

144    orthographic variations in collector names, assigning different text strings (i.e., recordedBy

145    strings in DarwinCore; hereafter "collector strings") to a single collector can be difficult and time

146    consuming for large datasets (Groom *et al.*, 2022). Thanks to the large-scale availability of

147    digitized historical and genealogical records (e.g., Ancestry.com, MyHeritage.com, and

148    Newspapers.com) and recent initiatives by historians of science to identify and disambiguate the

149    names of people who collected natural history specimens (e.g., Bionomia; Shorthouse, 2024;

150    Weeks *et al.*, 2024), we are for the first time able to identify and reconstruct what we call

151    oeuvres—all of the specimens a person has collected—of all contributors to a regional flora.

152

153    To disambiguate collector strings, we extracted the first collector in each collector string,

154    separating what we consider the principal collector (henceforth referred to as the collector) from

155    any associated collectors. Although associated collectors are crucial parts of any collection

156    team and deserve proper credit for their efforts, we focused our analysis on principal collectors

157    in this initial study. Our rationale is that the principal collector is usually responsible for recording

158    field notes and is likely to take on the major role of depositing the specimens in an herbarium

159    collection. We then separated the collector strings into words using the `unnest_tokens()`

160    function from `tidytext` (Silge & Robinson, 2016) and concatenated these words in

161    alphabetical order to standardize different transcriptions of the same text (e.g., "C. F. Parker", "C

162    F Parker", and "Parker, C. F." would all become "c,f,parker"). We then merged all records with

163    identical concatenated strings and manually validated each cluster—merging records with

164    different concatenated strings that represent the same collector—to ensure that each cluster

165  represented a single collector. We used biographical information from historical and
166  genealogical databases (e.g., Ancestry.com and Newspapers.com) and databases of natural
167  history collectors (i.e., Bionomia and Harvard Index of Botanists; Shorthouse, 2024; Harvard
168  University Herbaria, 2024) to reconstruct the oeuvres of collectors that collected under multiple
169  names, including their spouses' names. For instance, we identified "Mrs. C. S. Phelps" as Ora
170  Almira Phelps (née Parker) who collected under the names Mrs. Charles Sheppard Phelps,
171  Orra A. Phelps, Mrs. O. P. Phelps, and Orra Parker Phelps.
172
173  We excluded any collector strings that were ambiguous either because of obvious transcription
174  errors that could not be verified with a digital image of the specimen or had limited information.
175  To ensure that we were not conflating multiple collectors, we excluded records with only initials
176  (e.g., C.A.B.), only a surname (e.g., Boice), or only the initial of the first name and the surname
177  (e.g., C. Boice; 233,321 records removed,1,078,508 records remaining). We then removed
178  duplicate specimens (i.e., specimens collected by the same collector with the same specimen
179  number in DarwinCore's `recordNumber` field) so that each collection event is represented by a
180  single specimen (89,251 records removed). We did not remove any specimens without a
181  transcribed specimen number (i.e., those with "s.n.", "sn", or a blank `recordNumber` field) since
182  we could not confirm that that these specimens were duplicate collections. This resulted in our
183  final dataset of 989,257 specimens (Table S2).
184

185  ***Temporal Bias***
186  To investigate temporal trends in botanical collections, we calculated the number of specimens,
187  distinct species, sampling localities, and active collectors for each year during 1781–2024. We
188  then evaluated the relationship between these metrics and the oeuvre size of each collector on
189  a continuous scale from less-prolific (small oeuvres) to more-prolific collectors (large oeuvres).
190  We investigated seasonal variations in collection intensity by comparing the number of
191  specimens collected in each month and analyzed how this distribution changed with respect to
192  the oeuvre size of the collector who gathered the specimen.
193

194  ***Spatial Bias***
195  We quantified spatial bias by gridding the georeferenced specimens into 10-km grid squares
196  (hereafter localities) to help mitigate the effects of batch georeferencing and create equal-area
197  polygons for comparison (Franklin & Miller, 2009; Schmidt *et al.*, 2023). We calculated the
198  revisitation proportion for each collector as the number of specimens per unique collecting

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

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

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

202  collectors.

203

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

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

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

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

208  increasingly larger subsets of the data in 10,000 specimen increments (i.e., after arranging by

209  oeuvre size, we extracted the first 10,000 specimens, first 20,000 specimens, 30,000

210  specimens, etc.). We assessed how spatial bias from collectors with different oeuvre sizes

211  differs from two different null models: we randomly ordered specimens from our dataset to

212  determine if collections by more- or less-prolific collectors are more spatially clustered than the

213  overall specimens (randomized specimens); and simulated a new dataset by randomly sampling

214  from all localities in the northeastern US to determine how collections differ from spatially

215  random collections (simulated random sampling).

216

217  ***Taxonomic Bias***

218  To determine the relative representation of different taxa in herbarium collections, we calculated

219  collection depth as the average number of specimens per species in a given taxon in the

220  northeastern US (i.e., total specimens/unique species for each genus and family based on the

221  `acceptedScientificName` field from GBIF). We evaluate taxon size on a continuous scale,

222  whereby taxa with fewer species in the northeastern US are considered smaller and those with

223  more species are considered larger. Taxa with higher collection depths were considered better

224  represented in herbaria.

225

226  To assess how frequently collectors collect a species that they have already collected, we

227  calculated the proportion of species re-collected by each collector (i.e., total specimens/unique

228  species for each collector). Collectors who collected many specimens of the same species

229  would have a high re-collection proportion while those that collected only one specimen of each

230  species would have a re-collection proportion of one.

231

232    To investigate whether some taxa (i.e., species, genera, and families) were favored by

233    collectors over other taxa, we plotted the number of collections per taxon against the number of

234    collectors who collected each taxon. We fit a generalized additive model (GAM) to these points

235    to estimate how many collectors we expected to have collected each taxon based on the total

236    number of specimens of that taxon. Taxa that fell above this GAM curve were collected by more

237    people than expected (hereafter, favored taxa) and taxa that fell below the curve were collected

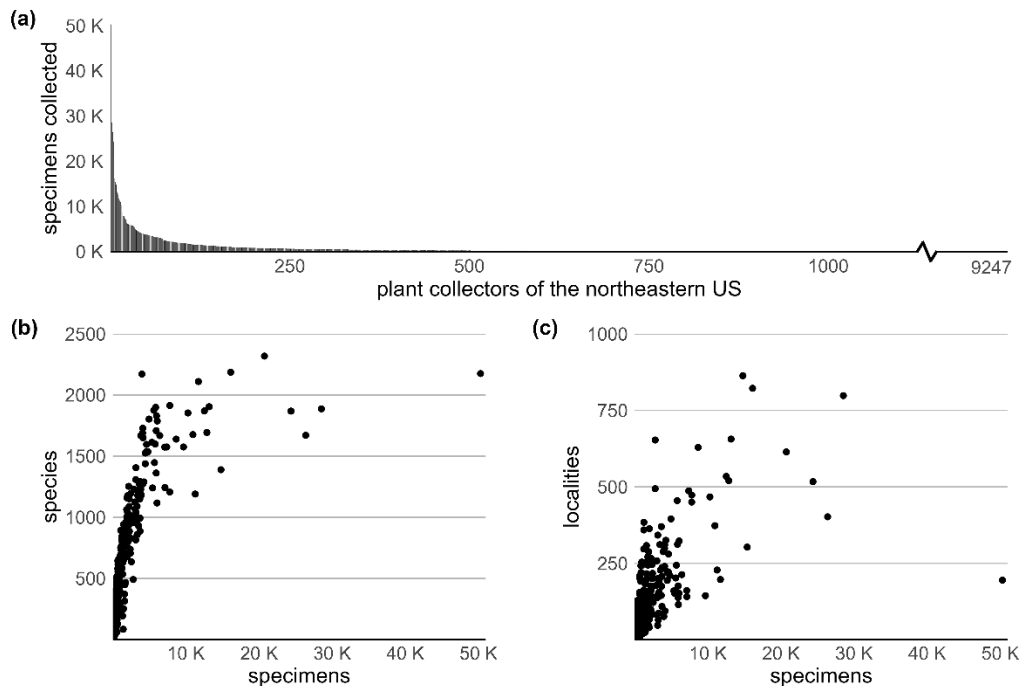238    by fewer people than expected (hereafter, commonly avoided taxa).

239

240    **Results**

241    *Collectors*

242    We identified 9247 collectors who collected plant specimens in the northeastern US (Fig. 1a;

243    Table S3). This is no doubt an underestimate of the total number of people who have

244    contributed to collections in the region since many collectors were excluded from our analysis

245    due to incomplete or ambiguous collector names and insufficient locality information (45% of

246    analyzable specimens removed) and those whose specimens have yet to be digitized. There

247    was a large variation in the number of specimens that each collector collected. We do not define

248    a threshold between more- and less-prolific collectors for any of our temporal, spatial, or

249    taxonomic analyses and instead evaluate variation in collector practices along a continuum of

250    oeuvre sizes with more- and less-prolific representing opposite ends of this range. However, we

251    briefly present results for some subsets of collectors below to demonstrate the overall variation

252    in the contributions of collectors with different oeuvre sizes. The vast majority (more than 90%;

253    8385 people) collected fewer than 100 specimens. Only 1.8% of collectors (171 people)

254    collected more than 1000 specimens (contributing 71% of the total number of collections). The

255    most prolific collector in our dataset was Robert L. Schaeffer, Jr., who collected 50,287

256    specimens (Fig. 1b). Half of all specimens from the northeastern US were collected by only 57

257    collectors (0.6% of collectors). Most collectors (70%; 6,549 people) collected fewer than ten

258    specimens (contributing 1.5% of collections).

259

260    People who collected less than 1000 specimens tended to collect only one specimen of each

261    species (Fig. 1b) and about ten specimens per locality (Fig. 1c). For collectors who collected

262    more than 1000 specimens, they tended to collect only one specimen for each species for the

263    first 1000 specimens they collected. After collecting about 1000 specimens, they collected

264    multiple specimens of the same species, and the number of species they collect plateaus near

265    2000 species. E. H. Eames collected the most plant species of any collector in our dataset

266 (2574 species, both vascular and nonvascular; Whelan, 1948). Most people collected either

267 vascular plants (85%) or non-vascular plants (7%), with only 8% collecting both types.
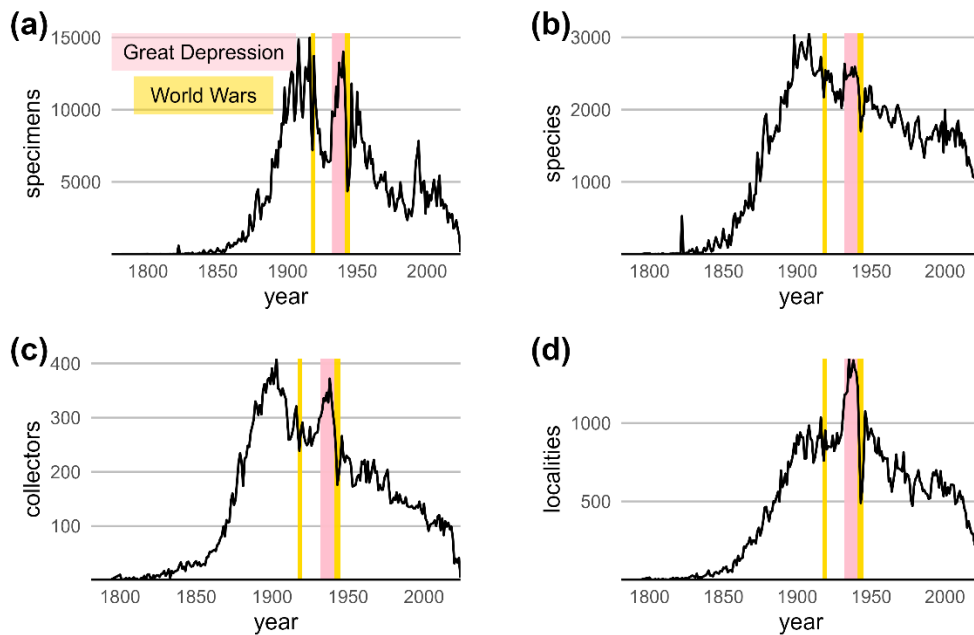
268



269

270 **Figure 1.** We identified 9247 people who collected herbarium specimens in the

271 northeastern US. The bar plot shows (a) the total number of unique specimens for each

272 plant collector in the northeastern US. The scatter plots show the relationship between

273 the number of specimens each person collected and (b) the number of species they

274 collected and (c) the number of localities in which they collected.
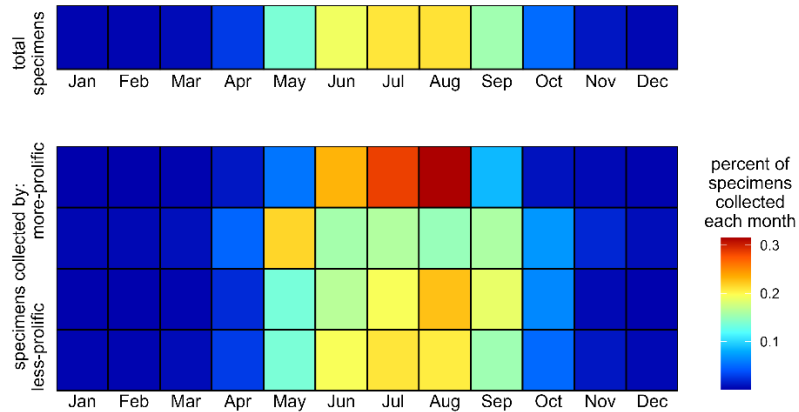
275

276 ***Temporal Bias***

277 The number of collectors active in a given year has varied substantially through time with peaks

278 during 1880–1916 and again during 1932–1941 (Fig. 2a–d). The number of active collectors is

279 strongly correlated with the number of specimens collected in a given year (cross-correlation

280 value of 0.90, p<0.001), species (0.94, p<0.001), and localities (0.90, p<0.001). The number of

281 specimens collected (Fig. 2a), species collected (Fig. 2b), and collectors active in a given year

282 (Fig. 2c) also peaked during 1880–1916 and 1935–1941 whereas the number of sampling

283 localities peaked only from 1935–1941 (Fig. 2d). All metrics have declined since 1950.

284

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



**Figure 2.** The line plots show the annual variation in (a) the number of specimens collected, (b) the number of species collected, (c) the number of active collectors, and (d) the number of localities in which specimens were collected from 1781–2024. The yellow bars indicate the years when the US was involved in World Wars I and II (1917–1920 and 1941–1946, respectively) and the pink bars represent the duration of environmental projects sponsored by the US federal government during the Great Depression (1929–1939).
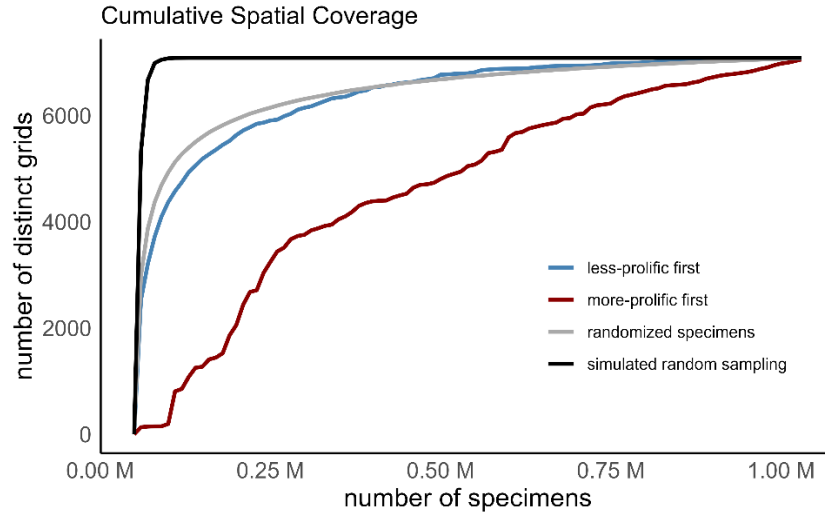
**Figure 3.** This graph shows the percentage of specimens collected in each month for all specimens (total specimens) and divided into quartiles based on oeuvre size (i.e., going from the first quartile of specimens collected by the least prolific collectors at the bottom to the fourth quartile of specimens collected by the most prolific collectors at the top).
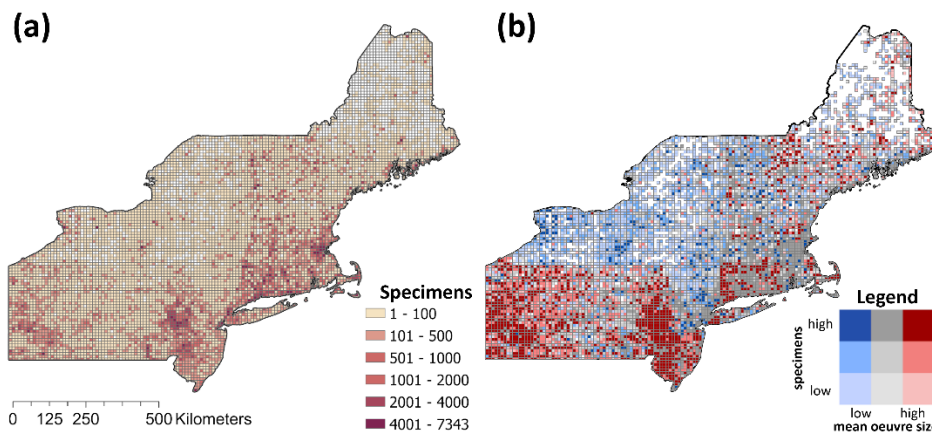
## *Spatial Bias*

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

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

Cumulative Spatial Coverage

**Figure 4.** Accumulation curves for the cumulative spatial coverage of gridded herbarium specimens based on the oeuvre size for collectors. Each curve contains all 989,257 specimens in our dataset with specimens added in different orders to demonstrate differences in the spatial coverage of specimens collected by more- and less-prolific 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).



**Figure 5.** The maps show (a) the density of collections in the northeastern US and (b) the relationship between collection density and areas where collections have been driven primarily by less-prolific collectors (blue; localities in the lowest tercile based on

12

337      mean oeuvre size), more-prolific collectors (red; localities in the middle tercile based on

338      mean oeuvre size ), or a mix of collector types (gray; localities in the top tercile based on
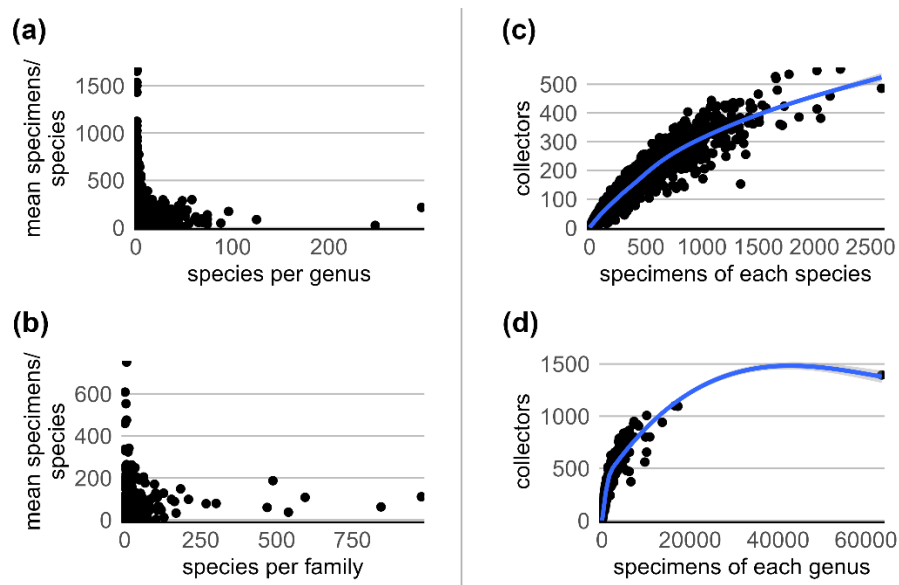
339      mean oeuvre size).

340

341    ***Taxonomic Bias***

342    Smaller genera are more likely to have a greater collection depth than larger genera; the same

343    is the case for smaller families (Fig. 6). Despite the overrepresentation of smaller genera,

344    several of the most frequently collected species are from large genera (e.g., three species of

345    *Carex*; for a list of the hundred most frequently collected species, see Table S4). Ferns are

346    dramatically overrepresented among the most frequently collected species (11 of the top 20

347    collected species were ferns). Within each year, 90% of specimens were collected during May–

348    September but only 46% of species were collected only during these five months. Species that

349    have been collected outside of the peak collection window (i.e., with at least one collection

350    during October–April) are far more likely to be overrepresented in herbaria compared with

351    species that have not been collected outside of peak collection months (Fig. S2). These non-

352    peak species include all but 18 of the 1000 most commonly collected species in the Northeast;

353    11 of these 18 are species of *Carex*. Despite also being collected in off-peak months, the top

354    species have been preferentially collected throughout the year, including during peak months;

355    96% of the top 1000 most collected species remain in the top 1000 when only collections from

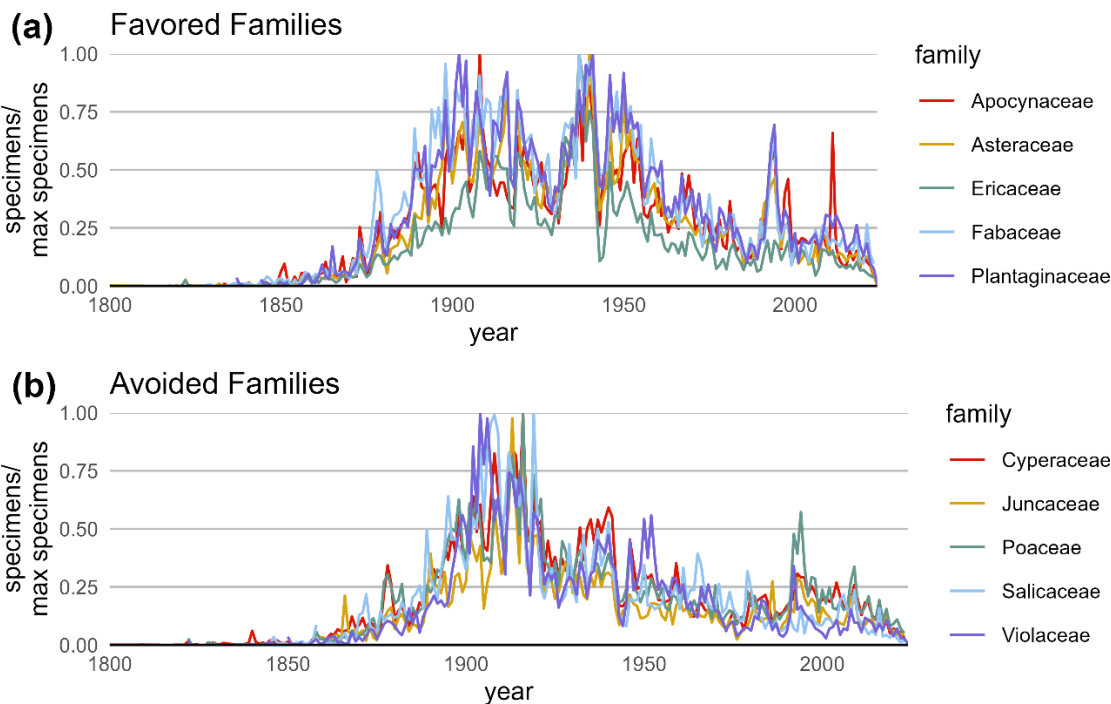356    peak months are considered.

357

358    Some species are overrepresented in collections because they were collected by many people

359    (e.g., *Arisaema triphyllum* (L.) Schott, *Onoclea sensibilis* L., and *Polystichum acrostichoides*

360    (Michx.) Schott; see Table S5 for information about the number of people who collected each

361    taxa mentioned in this section), whereas others are overrepresented because they were

362    collected intensively by a few people (e.g., *Sceptridium dissectum* (Spreng.) Lyon, *Scirpus*

363    *cyperinus* (L.) Kunth, and *Viola sororia* Willd.).

364

365    Some species were collected by far more people than expected from our GAM model (e.g.,

366    *Cypripedium acaule* Aiton and *Solanum dulcamara* L.) whereas *Dichanthelium acuminatum*

367    (Sw.) Gould & C.A.Clark was collected by far fewer people than expected. Similarly, some

368    genera were collected by more people than expected from our GAM model (e.g., *Lobelia,*

369    *Lysimachia*, and *Trifolium*), whereas others by fewer people than expected (e.g., *Crataegus,*

370    *Dichanthelium, Potamogeton, Salix, Sphagnum*). Some families were also collected by more

371    people than expected from our model (e.g., Apocynaceae, Asteraceae, Ericaceae, Fabaceae,

372    and Orchidaceae) and others by fewer than expected (e.g., Cyperaceae, Poaceae, Juncaceae,

373    Salicaceae, and Violaceae). Commonly favored families—collected by more people than

374    expected—typically had peaks in annual collections in the 1910s and 1930s, mirroring overall

375    trends in collections through time (Fig. 7). Commonly avoided families—collected by fewer

376    people than expected—typically had only a single peak during the 1910s. Some commonly

377    avoided families (e.g., Potamogetonaceae and Sphagnaceae), had relatively low collections

378    through time and its peaks correspond to specialist collectors rather than overall trends in

379    collections.

380



381

382    **Figure 6.** The plots show the collection depth (average number of specimens per

383    species) for each (a) genus and (b) family. The scatter plots in the right pane (panels c &

384    d) show the relationship between the number of specimens per species and the number

385    of collectors who collected these species of each (c) species and (d) genus.

386

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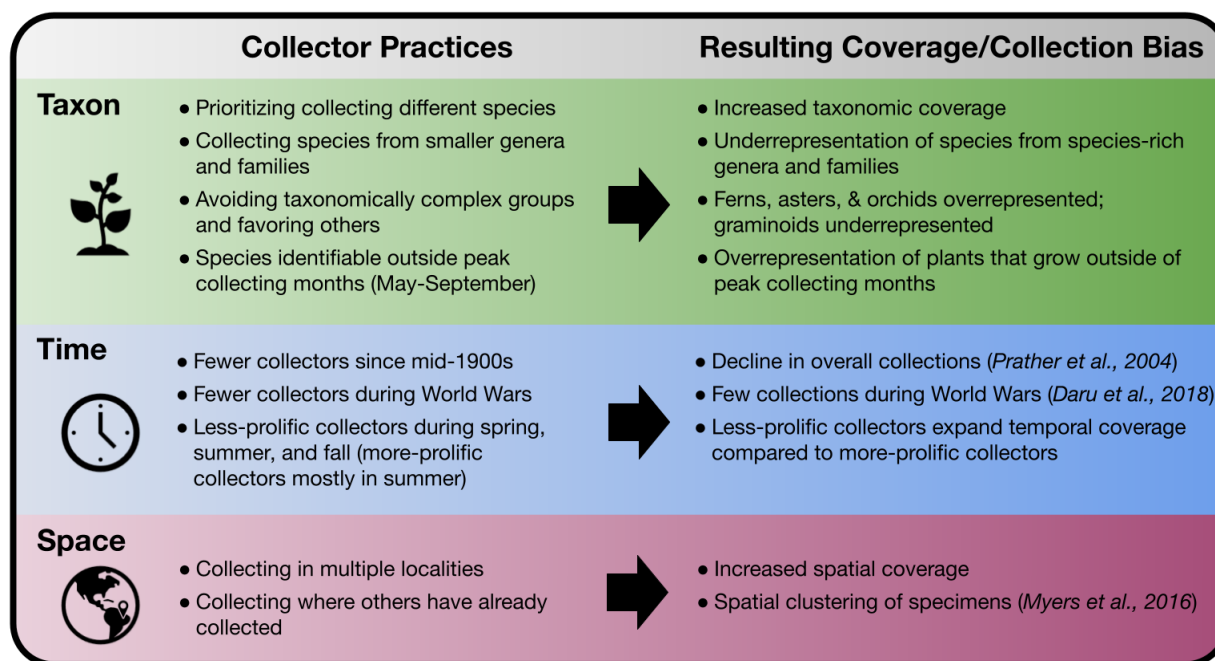
### *Summary of Results*

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

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

407    practices of all collectors: they tend to collect a.) more species rather than multiple specimens of

408    the same species; b.) about 10 specimens per locality during their lifetime; c.) from localities

409    sampled by other collectors; d.) during the peak growing season in spring and summer when

410    climates are more favorable and photosynthetic rates and reproduction are generally higher; e.)

411    species from smaller genera and families; and f.) particular species that are available outside of

412    peak collecting months (i.e., when climates are less favorable for plant growth. We also

413    identified that some collections norms have changed through time with collectors avoiding

414    several taxonomically complex taxa during some decades.

415

416    In contrast to the collections norms detailed above, we also identified several divergences

417    between the collector practices of more- versus less-prolific collectors. Specifically, more-prolific

418    collectors i.) collected largely during fewer months; ii.) had stronger affinities to certain localities;

419    and iii.) were not active in several large regions sampled by less-prolific collectors (e.g., the

420    state of New York, USA).

421

422    A summary of our findings is presented in Fig. 8, where we outline the collector practices and

423    resulting collection biases we have identified in the context of three key dimensions of bias:

424    taxon, time, and space. We include two previously identified temporal collection biases, the

425    decline in overall collections that was first presented by Prather *et al.* (2004) and the decline in

426    collections during World Wars I and II identified by Daru *et al.* (2018). We also include the

427    overall spatial clustering of collections, which was first defined by Myers *et al.* (2016).

428

| | Collector Practices | Resulting Coverage/Collection Bias |
|---|---|---|
| **Taxon** | • Prioritizing collecting different species<br>• Collecting species from smaller genera and families<br>• Avoiding taxonomically complex groups and favoring others<br>• Species identifiable outside peak collecting months (May-September) | • Increased taxonomic coverage<br>• Underrepresentation of species from species-rich genera and families<br>• Ferns, asters, & orchids overrepresented; graminoids underrepresented<br>• Overrepresentation of plants that grow outside of peak collecting months |
| **Time** | • Fewer collectors since mid-1900s<br>• Fewer collectors during World Wars<br>• Less-prolific collectors during spring, summer, and fall (more-prolific collectors mostly in summer) | • Decline in overall collections (*Prather et al., 2004*)<br>• Few collections during World Wars (*Daru et al., 2018*)<br>• Less-prolific collectors expand temporal coverage compared to more-prolific collectors |
| **Space** | • Collecting in multiple localities<br>• Collecting where others have already collected | • Increased spatial coverage<br>• Spatial clustering of specimens (*Myers et al., 2016*) |

429

430 **Figure 8.** This graphic describes the collector practices that have shaped overall the

431 overall coverage and collection bias in natural history collections along three

432 dimensions: taxon, time, and space.

433

434 **Discussion**

435 *Taxonomic bias: prioritizing greater species diversity and the underrepresentation of*

436 *large, complex taxa*

437 We found that botanists in the northeastern US prioritized collecting more species versus

438 collecting multiple specimens of the same species. Although this tendency has been viewed as

439 problematic in biology (Lewin, 1982; May, 2004), we assert that such collecting has contributed

440 considerably to expanding taxonomic coverage represented in NHCs and improving our

441 understanding of species diversity and distributions (Alba *et al.*, 2021). Despite this tendency,

442 however, collectors do not sample species randomly: many collect the same taxa while avoiding

443 others (Fig. 6). For instance, the brightly colored pink lady's slipper orchid (*Cypripedium acaule*

444 Aiton) was collected by many people (444 people) whereas hairy panicgrasses (*Dichanthelium*

445 *acuminatum* (Sw.) Gould & C.A.Clark) was collected by relatively few (153 people). This

446 collection norm affects our attempts to model biodiversity owing to the gap between taxon

447 diversity and abundance information recorded in NHCs versus their actual diversities and

448 abundances in nature (Elith & Leathwick, 2007; Gomes *et al.*, 2018). This pattern mirrors the

449 collection norm whereby collectors tend to collect ten specimens per locality and suggests that

450    collectors travelled to different localities to collect new species rather than comprehensively

451    collecting at a single locality.

452

453    Taxonomic collection norms have likely contributed to the overrepresentation of less species-

454    rich taxa with distinctive morphologies (e.g., *Lobelia, Polystichum,* and *Dryopteris*) in herbaria

455    relative to larger taxa that are often taxonomically challenging (e.g., *Carex, Crataegus,* and

456    *Salix*). Specimens from many large taxa were collected by fewer people than expected,

457    suggesting these were mainly collected by botanists with specialized taxonomic interests. In the

458    northeastern US, such specialist-prone taxa include genera like *Sphagnum* (peat mosses),

459    *Dichanthelium* (rosette grasses), *Salix* (willows), and *Crataegus* (hawthorns), and families like

460    Poaceae, Cyperaceae, and Juncaceae (collectively, the graminoids). These groups often

461    require microscopic examination to distinguish subtle differences necessary for accurate

462    species identification and often can only be identified with reproductive features at specific

463    maturation stages (FNA Editorial Committee, eds., 1993+). Further complicating species

464    identification and delimitation are their complex evolutionary histories, including infrageneric

465    hybridization (Ennos *et al.*, 2005). We hypothesize that this taxonomic bias in collections is often

466    driven by the perceived taxonomic complexity and difficulty to identify species within such

467    groups (for discussions of taxonomic complexity, see Ennos *et al.*, 2005; Karbstein *et al.*, 2024).

468    This collection norm suggests that the most diverse groups, which are likely in greatest need of

469    study, are woefully underrepresented in NHCs.

470

471    We also identified clear trends in shifting taxonomic collection norms through time, a pattern that

472    has received little attention. We observed that taxonomic biases have apparently shifted, with

473    certain taxa being favored and others apparently avoided across different generations of

474    botanists. For example, in the northeastern US, many collectors in the 1930s avoided families

475    like Poaceae, Cyperaceae, Juncaceae, and Sphagnaceae. We hypothesize that collectors from

476    the Citizens Conservation Corps, many of whom lacked formal botanical training, may have

477    avoided families they perceived as more complex. In other words, we hypothesize that

478    collectors are less prone to collect what they don't know. This has significant implications for

479    comparing temporal trends between taxa; variations in historical collection intensity may affect

480    apparent changes in characteristics such as species distribution modeling (Franklin & Miller,

481    2009) and phenology (Miller-Rushing *et al.*, 2008). Therefore, understanding the overall

482    temporal distribution of collections is crucial for appreciating how record availability—and the

483    uncertainty in these data—changes over time.

484

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

We identified an important divergent collection practice between more- and less-prolific collectors whereby less-prolific collectors contribute unique spatial coverage versus collections by more-prolific collectors (see Fig. 4). These less-prolific collectors enhance sampling near commonly collected localities (Fig. 4) and act as the backbone for entire regions where more-prolific collectors have not collected (Fig. 5b). For example, less-prolific collectors greatly improve spatial coverage in large portions of New York State, western Massachusetts, and near major universities (e.g., Rutgers University and Cornell University), likely highlighting the impact of student collectors on overall spatial sampling in the Northeast (see Fig. 5b). Thus, the cumulative spatial coverage by more-prolific collectors is considerably lower than that of less-prolific collectors, indicating that the collections made by the latter more accurately reflect plant diversity across different regions. It is important to note, however, that although there has been extensive herbarium digitization in the northeastern US, digitization is still ongoing and some of the spatial patterns apparent in currently available GBIF data will inevitably change as more data becomes available. In particular, several large collections in the region (e.g., the New York State Museum and the Pennsylvania State University Herbarium) do not publish their specimen data to major biodiversity aggregators like GBIF, potentially contributing to the comparably low specimen density in areas like upstate New York (Fig. 5a). Continued efforts to digitize 'silent' herbaria (Zhigila, Schmidt *et al.*, 2025) and make their data digitally accessible are necessary for understanding how data from NHCs can be leveraged for studying global biodiversity. Interestingly, the spatial bias of less-prolific collectors does not differ significantly from the overall spatial bias in herbaria. However, these collections are still biased with respect to random sampling. This suggests that while less-prolific collectors do not exhibit the same preference for specific collection sites as more-prolific collectors, they also tend to revisit locations where collections have previously been made. Despite this spatial collection norm, the increased spatial coverage provided by less-prolific collectors has greatly improved the overall spatial sampling in herbaria. This increased spatial coverage has helped facilitate the recent application of herbarium data to disciplines that rely on extensive sampling; for example, ecology (Meineke *et al.*, 2019a; Heberling, 2022); invasion biology (Crawford & Hoagland 2009; Schmidt *et al.,* 2023), species distribution modeling (Daru *et al.*, 2021), environmental science (Carbone *et al.*, 2023; Jakovljević *et al.*, 2024), and conservation biology (Schatz, 2002).

517

518    Finally, the broad spatial sampling by numerous less-prolific collectors that we identified reflects

519    patterns also observed with contemporary iNaturalist data, where contributions by millions of

520    community scientists greatly extend spatial sampling beyond what is captured in herbaria

521    (Eckert *et al.*, 2024; also see Daru & Rodriguez, 2023). This similarity indicates that the spatial

522    biases of community scientists align more closely with those of less-prolific collectors than with

523    the more-prolific collectors who contributed heavily to overall spatial biases in collections.

524    Similarly, several studies have demonstrated that small, regional collections provide unique

525    temporal and spatial coverage not represented in larger collections (Glon *et al.*, 2017; Monfils *et*

526    *al.*, 2020; Marsico *et al.*, 2020). We hypothesize that the expanded coverage of smaller herbaria

527    is driven by the efforts of less-prolific collectors who also provide unique temporal and spatial

528    coverage not captured by more-prolific collectors.

529

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

531    The substantial declines in collections over the past 75 years is consistent with trends observed

532    in other regional floras (Prather *et al.*, 2004; Daru *et al.*, 2018) and is strongly correlated with

533    declines in the number of active collectors. This suggests that while more-prolific collectors may

534    heavily influence the interannual intensity of collections at certain times (Bebber *et al.*, 2012;

535    Daru *et al.*, 2018), the overall trends are primarily driven by fluctuations in the number of all

536    active collectors.

537

538    Notably, the reduction in annual collections coincided with the years when the US was involved

539    in World Wars I (1917–1920) and II (1941–1946). During the two world wars, millions of men

540    were conscripted for military service and at the same time millions of women, students, and

541    older Americans entered the workforce (Witt, 1942; Wilcock, 1957) and would have been largely

542    unable to collect plants.. Following decreased collections during World War I, the spike in

543    collections and active collectors from 1932 through 1941 corresponds with US government

544    efforts to reduce unemployment and support environmental projects during the Great

545    Depression (1929–1939; Salmond, 1967). During this period, the government employed

546    thousands of citizens—primarily young men aged 18 to 25—for projects focusing on

547    environmental improvements (e.g., in the Civilian Conservation Corps; Salmond, 1967). A key

548    objective of these initiatives was to produce local species inventories, documented through

549    "complete herbaria," to aid in land planning and protection (Department of the Interior, 1936).

550    Since these projects often targeted similar habitats—primarily forested areas—many inventories

551    likely covered areas with similar species composition in the northeastern US. Consequently,

552 despite the spikes in collections, active collectors, and collection locations during this time, the
553 number of species collected during this period did not increase substantially. Once World War II
554 began and people from the same demographic were heavily drafted into WWII, all metrics once
555 again quickly declined. This highlights how major socio-political events affecting significant
556 population segments can directly impact NHCs by reducing the pool of available collectors.
557 Similar impacts of socio-political events on NHCs were recently documented in collection
558 requests for multiomic sampling, which plummeted during the global COVID pandemic (Davis *et*
559 *al.*, 2024).
560
561 We identified that less-prolific collectors increased overall sampling at the start and end of the
562 primary growing season (late spring and early autumn), which diverges from collections by
563 more-prolific collectors whose activity during these periods markedly decreases. The intensity of
564 sampling during these off-peak periods is crucial for improving the accuracy of phenological
565 estimates (Miller-Rushing *et al.*, 2008) and understanding the impact of anthropogenic climate
566 change on early- and late-season species (Kudo & Ida, 2013; Park *et al.*, 2023). We
567 hypothesize that the increased sampling by less-prolific collectors at the beginning and end of
568 the growing season (i.e., April–May and September–October) might be related to student
569 collections in university botany classes during the academic year (typically September–May).
570
571 Surprisingly, although 90% of specimens are collected in the northeastern US between May and
572 September, species collected outside the peak months are disproportionately represented
573 among the most abundant species in herbaria. These include many evergreen (e.g.,
574 *Polystichum acrostichoides* (Michx.) Schott and *Dryopteris marginalis* (L.) A.Gray), woody (e.g.,
575 *Vaccinium corymbosum* L. and *Acer rubrum* L.), and early-flowering species (e.g., *Viola sororia*
576 Willd. and *Arisaema triphyllum* (L.) Schott), as well as species with winter-available flowers or
577 fruits (e.g., *Ilex verticillata* (L.) A.Gray and *Hamamelis virginiana* L.). We hypothesize this
578 overrepresentation is driven by collectors' familiarity with these species, which are more
579 accessible and—in some cases—more identifiable outside of peak collection months when
580 fewer species are available.
581

582 ***Exceptions to the norms: unique collector practices contribute overall bias***
583 Despite the similar collector practices we identified, we emphasize that understanding how
584 some collectors diverged from these norms is important for understanding overall collection bias
585 in NHCs. For example, the most prolific collector in our dataset, R. L. Schaeffer, Jr., collected

586    50,287 specimens from only 195 localities—far fewer than expected based on our model. He

587    collected, almost exclusively, in the vicinity of Allentown, PA where Schaeffer taught botany at

588    Muhlenberg College from 1954-1983 ('R. L. Schaeffer Obituary', 2001). His singular efforts had

589    an outsized impact on overall spatial bias in the northeastern US with his collections being the

590    main driver of the high collection density in eastern PA, one of the most collection-dense areas

591    in the northeastern US. Furthermore, the expansive taxonomic coverage and high collection

592    depth of Schaeffer's specimens provides a rich documentation of the flora of eastern

593    Pennsylvania over nearly a half century that can be leveraged for a diversity of collections-

594    based investigations (e.g., Meineke *et al.*, 2019b). This highlights how integrating historical

595    information about collectors (especially mega-collectors like Schaeffer) can help explain the

596    more stochastic processes in biodiversity data and can illuminate important datasets better

597    characterizing species and ecosystem responses to anthropogenic pressures.

598

599    *Conclusion*

600    Our findings reveal how our understanding of biodiversity is founded on the cumulative effort of

601    thousands of people, many of whom have made small but impactful contributions to natural

602    history collections (NHCs). The cumulative spatial, temporal, and taxonomic practices of all

603    collectors give rise to the overall biases in collections. It is crucial that we identify and categorize

604    these collector practices to better understand the drivers of overall collection bias in NHCs and

605    begin developing tools to address them. We have identified numerous predictable collection

606    norms that appear to have shaped overall bias in NHCs. The predictability of these biases

607    provides an exciting and promising opportunity to begin incorporating statistical tools to address

608    collection biases in biodiversity models. These results can also be leveraged to guide future

609    collection efforts that can minimize gaps in collections and reduce bias in NHCs moving forward.

610    We highlight that collector practices—even by those who collected only a small number of

611    specimens—have vastly expanded the coverage of NHCs and we assert that continued

612    collections of all sizes are crucial for continuing to expand the coverage of NHCs and further

613    increasing their utility for understanding biodiversity in the face of global change.

614

615    **<u>Acknowledgements</u>**

22

**Author Contribution**

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

**Data Availability Statement**

The data generated during this study are available in the supporting information of this
manuscript. Tables S2, S3, and all code created for this study are available on the Harvard
Dataverse (https://doi.org/10.7910/DVN/OJCODH).

**Conflict of Interest Statement**

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

**References**

**Alba C, Levy R, Hufft R**. **2021**. Combining botanical collections and ecological data to better
describe plant community diversity (J-Z Wan, Ed.). *PLOS ONE* **16**: e0244982.

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

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

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

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

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

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

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

**Daru BH, Rodriguez J**. **2023**. Mass production of unvouchered records fails to represent global biodiversity patterns. *Nature Ecology & Evolution* **7**: 816–831.

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

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

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

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

**Eckert I, Bruneau A, Metsger DA, Joly S, Dickinson TA, Pollock LJ**. **2024**. Herbarium collections remain essential in the age of community science. *Nature Communications* **15**: 7586.

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

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

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

683 **Flora of North America Editorial Committee, eds. 1993**. Flora of North America (FNA). *Flora*
684 *of North America*.

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

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

689 **GBIF.org. 2024.** GBIF Occurrence Download. Available at: https://doi.org/10.15468/dl.rndw9f.
690 (Accessed: 28 August 2024).

691 **Glon HE, Heumann BW, Carter JR, Bartek JM, Monfils AK**. **2017**. The contribution of small
692 collections to species distribution modelling: A case study from Fuireneae (Cyperaceae).
693 *Ecological Informatics* **42**: 67–78.

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

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

701 **Harvard University Herbaria**. **2024**. Harvard Index of Botanists. *Harvard University Herbaria &*
702 *Libraries*.

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

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

708 **Jakovljević K, Mišljenović T, Van Der Ent A, Baker AJM, Invernón VR, Echevarria G**. **2024**.

709 "Mining" the herbarium for hyperaccumulators: Discoveries of nickel and zinc

710 (hyper)accumulation in the genus *NOCCAEA* (Brassicaceae) through X-ray fluorescence

711 herbarium scanning. *Ecological Research* **39**: 450–459.

712 **Johnson KR, Owens IFP, the Global Collection Group**. **2023**. A global approach for natural

713 history museum collections. *Science* **379**: 1192–1194.

714 **Karbstein K, Kösters L, Hodač L, Hofmann M, Hörandl E, Tomasello S, Wagner ND,**

715 **Emerson BC, Albach DC, Scheu S, *et al.* 2024**. Species delimitation 4.0: integrative taxonomy

716 meets artificial intelligence. *Trends in Ecology & Evolution* **39**: 771–784.

717 **Kozlov MV, Sokolova IV, Zverev V, Zvereva EL**. **2021**. Changes in plant collection practices

718 from the 16th to 21st centuries: implications for the use of herbarium specimens in global

719 change research. *Annals of Botany* **127**: 865–873.

720 **Kudo G, Ida TY**. **2013**. Early onset of spring increases the phenological mismatch between

721 plants and pollinators. *Ecology* **94**: 2311–2320.

722 **Lendemer J, Thiers B, Monfils AK, Zaspel J, Ellwood ER, Bentley A, LeVan K, Bates J,**

723 **Jennings D, Contreras D, *et al.* 2020**. The Extended Specimen Network: A Strategy to

724 Enhance US Biodiversity Collections, Promote Research and Education. *BioScience* **70**: 23–30.

725 **Lewin R**. **1982**. Biology Is Not Postage Stamp Collecting: Ernst Mayr, the eminent Harvard

726 evolutionist, explains why he thinks some physical scientists have a problem with evolution.

727 *Science* **216**: 718–720.

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

729 guidelines.

730 **Marín-Rodulfo M, Rondinel-Mendoza KV, Martín-Girela I, Cañadas EM, Lorite J**. **2024**. Old

731 meets new: Innovative and evolving uses of herbaria over time as revealed by a literature

732 review. *PLANTS, PEOPLE, PLANET* **6**: 1261–1271.

733 **Marsico TD, Krimmel ER, Carter JR, Gillespie EL, Lowe PD, McCauley R, Morris AB,**

734 **Nelson G, Smith M, Soteropoulos DL, *et al.* 2020**. Small herbaria contribute unique

735 biogeographic records to county, locality, and temporal scales. *American Journal of Botany* **107**:

736 1577–1587.

737 **May RM**. **2004**. Tomorrow's taxonomy: collecting new species in the field will remain the rate–

738 limiting step (HCJ Godfray and S Knapp, Eds.). *Philosophical Transactions of the Royal Society*

739 *of London. Series B: Biological Sciences* **359**: 733–734.

740 **Meineke EK, Classen AT, Sanders NJ, Jonathan Davies T**. **2019a**. Herbarium specimens

741 reveal increasing herbivory over the past century (A Iler, Ed.). *Journal of Ecology* **107**: 105–117.

742 **Meineke EK, Davies TJ, Daru BH, Davis CC**. **2019b**. Biological collections for understanding

743 biodiversity in the Anthropocene. *Philosophical Transactions of the Royal Society B: Biological*

744 *Sciences* **374**: 20170386.

745 **Miller-Rushing AJ, Inouye DW, Primack RB**. **2008**. How well do first flowering dates measure

746 plant responses to climate change? The effects of population size and sampling frequency.

747 *Journal of Ecology* **96**: 1289–1296.

748 **Monfils AK, Krimmel ER, Bates JM, Bauer JE, Belitz MW, Cahill BC, Caywood AM, Cobb**

749 **NS, Colby JB, Ellis SA,** *et al.* **2020**. Regional Collections Are an Essential Component of

750 Biodiversity Research Infrastructure. *BioScience* **70**: 1045–1047.

751 **Myers N, Mittermeier RA, Mittermeier CG, Da Fonseca GAB, Kent J. 2000. Biodiversity**

752 **hotspots for conservation priorities. *Nature* 403: 853–858.**

753 **Park DS, Davis CC**. **2017**. Implications and alternatives of assigning climate data to

754 geographical centroids. *Journal of Biogeography* **44**: 2188–2198.

755 **Park DS, Xie Y, Ellison AM, Lyra GM, Davis CC**. **2023**. Complex climate-mediated effects of

756 urbanization on plant reproductive phenology and frost risk. *New Phytologist* **239**: 2153–2165.

757 **Pebesma E**. **2018**. Simple Features for R: Standardized Support for Spatial Vector Data. *The R*

758 *Journal* **10**: 439.

759 **Pebesma E, Bivand R**. **2023**. *Spatial Data Science: With Applications in R*. New York:

760 Chapman and Hall/CRC.

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

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

763 Sons.

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

767 **Preston CD**. **2013**. Following the BSBI's lead: the influence of the *Atlas of the British flora* ,
768 1962–2012. *New Journal of Botany* **3**: 2–14.

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

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

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

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

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

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

782 **Shorthouse DP**. **2024**. Bionomia. *Bionomia*.

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

785 **Sweeney PW, Starly B, Morris PJ, Xu Y, Jones A, Radhakrishnan S, Grassa CJ, Davis CC**.
786 **2018**. Large–scale digitization of herbarium specimens: Development and usage of an
787 automated, high–throughput conveyor system. *TAXON* **67**: 165–178.

788 **United States Census Bureau. 2024.** Cartographic Boundary Files: States: 1 : 500,000
789 (national). Available at: https://www.census.gov/geographies/mapping-files/time-
790 series/geo/cartographic-boundary.html. (Accessed: 6 September 2024).
791

792    **Webster MS (Ed.)**. **2017**. *The extended specimen: emerging frontiers in collections-based*

793    *Ornithological Research*. Boca Raton London New York: CRC Press, Taylor & Francis Group.

794    **Weeks A, Collins E, Majors T, Murrell Z, Paul D, Sheik M, Shorthouse D, Zeringue-**

795    **Krosnick S**. **2024**. Workshop Report: Supporting inclusive and sustainable collections-based

796    research infrastructure for systematics (SISRIS). *Research Ideas and Outcomes* **10**: e126532.

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

798    **Wilcock RC**. **1957**. The Secondary Labor Force and the Measurement of Unemployment. In:

799    The Measurement and Behavior of Unemployment. National Bureau for Economic Research,

800    167–210.

801    **Witt B**. **1942**. Labor in Transition to a War Economy. In: Hanna HS, ed. Monthly Labor Review.

802    Bureau of Labor Statistics.

803    **Zhigila D, Schmidt RJ, Thiers B, Abdul S, Abdullahi S, AbdulRahaman A, Aigbokhan E,**

804    **Ajibade G, Ajikah L, Akomaye F,** *et al.* **2025**. Biodiversity science is improved when silent

805    herbaria speak.

806

807    **Supporting Information**

808    **Table S1** Herbaria whose specimens were used for this study, indicating the institution code,

809    institution name, and the number of specimens from each herbarium that were used in this

810    study.

811    **Table S2** Total specimens used in this study after data cleaning, georeferencing, and collector

812    disambiguation.

813    **Table S3** A table containing the DarwinCore recordedBy strings from gbif, the unique identifier

814    representing each collector, and the number of specimens, species, and localities in which each

815    person collected plants.

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

817    **Table S5** The number of specimens, species, and collectors that collected taxa mentioned in

818    the text of the manuscript.

819    **Fig. S1** A flowchart showing the data cleaning process including the number of specimens

820    removed at each step.

821    **Fig. S2** A boxplot showing the difference in number of specimens of each species related to

822    whether the species has been collected only during peak collection months (May, June, July,

823    August, and September) or also collected in non-peak months.