

1 **The Neotoma Paleoecology Database: A multi-proxy, international**
2 **community-curated data resource**

3

4 **Author Names & Affiliations**

5 Williams, John W.^{1,2,*}
6 Grimm, Eric C.^{3,*}
7 Blois, Jessica⁴
8 Charles, Donald F.⁵
9 Davis, Edward⁶
10 Goring, Simon J.¹
11 Graham, Russell W.⁷
12 Smith, Alison J.⁸
13 Anderson, Michael⁹
14 Arroyo-Cabralles, Joaquin¹⁰
15 Ashworth, Allan C.¹¹
16 Betancourt, Julio L.¹²
17 Bills, Brian W.¹³
18 Booth, Robert K.¹⁴
19 Buckland, Philip¹⁵
20 Curry, B. Brandon¹⁶
21 Giesecke, Thomas¹⁷
22 Jackson, Stephen T.¹⁸
23 Latorre, Claudio¹⁹
24 Nichols, Jonathan²⁰
25 Purdum, Timshel²¹
26 Roth, Robert E.^{1, 22}
27 Stryker, Michael¹²
28 Takahara, Hikaru²³

29

30 ¹Department of Geography, University of Wisconsin-Madison, Madison, WI 53706

31 ²Center for Climatic Research, University of Wisconsin-Madison, Madison, WI 53706

32 ³Department of Earth Sciences, University of Minnesota, Minneapolis, MN 55455

33 ⁴School of Natural Sciences, University of California, Merced, CA 95343

34 ⁵Earth and Environmental Science, Drexel University and Patrick Center, Academy of Natural
35 Sciences of Drexel University, Philadelphia, PA 19103

36 ⁶Department of Earth Sciences and Museum of Natural and Cultural History, University of
37 Oregon, Eugene, OR 97403

38 ⁷Department of Geosciences, College of Earth and Mineral Sciences, The Pennsylvania State
39 University, State College, PA 16802

40 ⁸Department of Geology, Kent State University, Kent, OH 44242

41 ⁹SpatialIT, State College, PA 16802

42 ¹⁰Laboratorio de Arqueozoología, Instituto Nacional de Antropología e Historia, 06060 Ciudad
43 de Mexico, CdMx

44 ¹¹Department of Geosciences, North Dakota State University, Fargo, ND 58108

45 ¹²National Research Program, Water Mission Area, U.S. Geological Survey, Reston VA 20192

46 ¹³Center for Environmental Informatics, The Pennsylvania State University, State College, PA
47 16802

48 ¹⁴Earth and Environmental Sciences Department, Lehigh University, Bethlehem, PA 18015

49 ¹⁵Environmental Archaeology Lab, Dept. of Historical, Philosophical & Religious Studies, Umeå
50 University, Umeå, SE-90187

51 ¹⁶Illinois State Geological Survey, Champaign, IL 61820

52 ¹⁷Department of Palynology and Climate Dynamics, Albrecht-von-Haller-Institute for Plant
53 Sciences, University of Göttingen, Göttingen

54 ¹⁸Southwest Climate Science Center, U.S. Geological Survey, Tucson, AZ 85721 and
55 Department of Geosciences, University of Arizona, Tucson, AZ 85721

56 ¹⁹Departamento de Ecología, Facultad de Ciencias Biológicas, Pontificia Universidad Católica de
57 Chile, Casilla 114-D, Santiago and Institute of Ecology & Biodiversity (IEB), Santiago, Chile

58 ²⁰Lamont-Doherty Earth Observatory, Palisades, NY 10964

59 ²¹Academy of Natural Sciences of Drexel University, Philadelphia, PA 19103

60 ²²Cartography Lab, University of Wisconsin-Madison, Madison, WI 53706

61 ²³Laboratory of Forest Vegetation Dynamics, Kyoto Prefectural University, Hangi-cho,
62 Shimogamo, Sakyo-ku, Kyoto 606-8522

63 *These authors contributed equally to this manuscript.

64

65 **Corresponding Author**

66 John W. Williams

67 jww@geography.wisc.edu

68 Department of Geography
69 550 North Park St
70 University of Wisconsin-Madison
71 Madison, WI 53706
72 608-265-5537 office
73 608-265-3994 fax
74

75 **ABSTRACT**

76 The Neotoma Paleoecology Database is a community-curated data resource that supports
77 interdisciplinary global change research by enabling broad-scale studies of taxon and
78 community diversity, distributions, and dynamics during the large environmental changes of the
79 past. By consolidating many kinds of data into a common repository, Neotoma lowers costs of
80 paleodata management, makes paleoecological data openly available, and offers a high-quality,
81 curated resource. Neotoma's distributed scientific governance model is flexible and scalable,
82 with many open pathways for participation by new members, data contributors, stewards, and
83 research communities. The Neotoma data model supports, or can be extended to support, any
84 kind of paleoecological or paleoenvironmental data from sedimentary archives. Data additions
85 to Neotoma are growing and now include >3.8 million observations, >17,000 datasets, and
86 >9,200 sites. Dataset types currently include fossil pollen, vertebrates, diatoms, ostracodes,
87 macroinvertebrates, plant macrofossils, insects, testate amoebae, geochronological data, and
88 the recently added organic biomarkers, stable isotopes, and specimen-level data. Multiple
89 avenues exist to obtain Neotoma data, including the Explorer map-based interface, an
90 Application Programming Interface, the *neotoma* R package, and digital object identifiers. As
91 the volume and variety of scientific data grow, community-curated data resources such as
92 Neotoma have become foundational infrastructure for big data science.

93 **KEYWORDS**

94 Biogeography; Geoinformatics; Global; Micropaleontology; Paleoclimatology; Paleodatabases;
95 Paleoecology; Paleoecoinformatics; Paleolimnology; Paleontology

96

97

99 **INTRODUCTION**

100 The Neotoma Paleoecology Database (hereafter called Neotoma, www.neotomadb.org)
101 was launched in 2009 with a mission to provide an open, community-curated, sustainable, and
102 high-quality repository for multiple kinds of paleoecological data. Neotoma's name refers to the
103 behavior of woodrats or packrats (genus *Neotoma*), which (inadvertently) serve paleoecology by
104 gathering diverse biological materials into their nests, there to be preserved for future
105 generations. Although Neotoma itself is relatively young, it builds upon decades of effort by
106 paleoecologists, paleoclimatologists, and paleontologists to gather individual records into larger
107 spatial networks for the purpose of studying ecological, evolutionary, biogeographic, climatic,
108 and cultural processes at spatial scales beyond the scope of any single site-level
109 paleoecological record. The gathering of these records is expensive, with substantial
110 investments in money and time, and their scientific value is multiplied when aggregated into
111 larger networks. By bringing these resources into a single open data resource with an
112 accompanying distributed governance framework, Neotoma seeks to accelerate our capacity to
113 do global-scale paleoscience, serve as an open-source platform for new kinds of analytics and
114 visualizations, enhance reproducibility, and increase the longevity and sustainability of our
115 communities' hard-won data.

116 Neotoma's creation and design is motivated by the goal of enabling global-scale science
117 from long-term, site-level data. Ecological processes operate across a wide range of interacting
118 spatial and temporal scales (Heffernan et al., 2014). Dynamics at one location are often
119 interpretable only in the context of larger-scale biogeographic and climatic processes (Webb,
120 1997), and ecosystems can be affected by slow processes that were triggered by events
121 centuries or even millennia ago (Goring and Williams, 2017; Svenning and Sandel, 2013).
122 Networks of paleoecological records, therefore, provide fundamental scientific infrastructure for
123 understanding the responses of species to large and abrupt environmental changes, the
124 mechanisms that promote resilience, and the interplay between climatic and biotic interactions
125 (Blois et al., 2013; Dawson et al., 2011; Jackson and Blois, 2015; Moritz and Agudo, 2013).
126 Examples include the processes controlling contemporary and past patterns of community,
127 species, and genetic diversity (Blarquez et al., 2014; Cinget et al., 2015; De La Torre et al.,
128 2014; Fritz et al., 2013; Gutiérrez-García et al., 2014; Jezkova et al., 2015; Sandom et al.,
129 2014); identification of species refugia (Bennett and Provan, 2008; Gavin et al., 2014; Vickers

130 and Buckland, 2015); rates of species expansion (Giesecke et al., 2017; Ordóñez and Williams,
131 2013); the reshuffling of species into no-analog communities during climate change (Finsinger
132 et al., 2017; Graham et al., 1996; Radeloff et al., 2015); the timing and patterns of abrupt
133 ecological and climate change (Seddon et al., 2015; Shuman et al., 2009), quantification of the
134 time lags between abrupt climate change and local ecological response (Ammann et al., 2013;
135 Birks, 2015); and the timing, causes, and consequences of late-Quaternary megafaunal
136 extinctions (Doughty et al., 2013; Emery-Wetherell et al., in press; Lorenzen et al., 2011).

137 The scientific communities interested in paleoecological data extend well beyond
138 paleoecology and biogeography. Paleoecological data such as fossil pollen, diatoms, and
139 marine foraminifera are the backbone of continental- to global-scale paleoclimatic
140 reconstructions developed to benchmark climate models and assess feedbacks within the earth
141 system (Bartlein et al., 2011; CLIMAP Project Members, 1976; Marcott et al., 2013; MARGO
142 Project Members, 2009; Shakun et al., 2012; Trouet et al., 2013; Viau et al., 2012; Wright et al.,
143 1993) and constrain estimates of climate sensitivity (Schmittner et al., 2011). Paleoecological
144 data help establish ecosystem baselines and trajectories for managers seeking to conserve
145 species and ecosystems of concern (Barnosky et al., 2017; Clarke and Lynch, 2016; Dietl et al.,
146 2015; Panagiotakopulu and Buchan, 2015; Whitehouse et al., 2008). Similarly, paleoecological
147 data are necessary for understanding the interactions between past environmental change and
148 early human evolution, land use, cultural and technological innovation, and dispersal at local to
149 global scales (deMenocal, 2001; Ellis et al., 2013; Gaillard et al., 2010; Grant et al., 2014;
150 Kaplan et al., 2011; Kaplan et al., 2009; Muñoz et al., 2010).

151 In response to these scientific drivers, multiple databases have been developed by
152 multiple teams over the past 30 years for different kinds of Pliocene-Quaternary fossil data.
153 These prior efforts, beginning in North America and Europe in the 1970s (Bernabo and Webb,
154 1977; Davis, 1976; Grimm et al., 2013; Huntley and Birks, 1983; Sadler et al., 1992), resulted in
155 multiple paleoecological databases, each usually restricted to a particular dataset type or
156 region, e.g. the North American, European, African, and Latin American Pollen Databases
157 (Flantua et al., 2015; Fyfe et al., 2009; Grimm et al., 2013; Vincens et al., 2007); the FAUNMAP
158 and MIOMAP terrestrial vertebrate databases (Carrasco et al., 2007; Graham et al., 1996), the
159 Mexican Quaternary Mammal Database (MQMD, Arroyo-Cabrales et al., 2007, 2009), NANODe
160 (Forester et al., 2005), the Diatom Paleolimnology Data Cooperative
161 (<https://diatom.anasp.org/dpdc/>) (Sullivan and Charles, 1994), the North American packrat
162 midden (Betancourt et al., 1990) and plant macrofossil databases (Jackson et al., 2000;
163 Jackson et al., 1997), the Base de Données Polliniques et Macrofossiles du Québec (Richard,

164 1995), the BUGS insect database (Sadler et al., 1992), and others. Other, newer proxies, such
165 as testate amoebae and organic biomarkers, are just beginning to be gathered for use in broad-
166 scale studies and need a common platform for data archival, sharing, and re-use.

167 The importance of these databases cannot be overstated. Hundreds of scientific papers
168 have utilized them (Fig. 1); and entire dissertations and subsequent papers have been based on
169 them (e.g. Buckland, 2007; Li, 2004; Lyons, 2001). Nevertheless, the various paleoecology
170 databases were typically established with either one-time or sporadically funded projects. Long-
171 term maintenance, sustainability, and development have plagued virtually all paleoecological
172 and paleoclimatic database efforts mainly because continuous funding is needed for both
173 information technology (IT) and data preparation and cleaning. Some databases have had
174 curatorial support from museums and government entities (e.g. Canada Museum of Nature's
175 curation of Delorme ostracode database), others have not. Funding hiatuses rarely cause these
176 databases to disappear entirely, but have caused data backlogs and delays, with long delays
177 between data contributions by individual scientists and their release to the public. Additionally,
178 early databases were maintained in flat files or standalone, desktop-database management
179 systems such as Paradox or Microsoft Access (Grimm et al., 2013), versus the new standard of
180 client/server systems that serve data over the internet. Copies of databases sometimes
181 proliferated and individual copies rapidly became obsolete (e.g. the initial release of FAUNMAP
182 was distributed by hard copy and floppy disks, as well as on-line, FAUNMAP Working Group,
183 1994). Because paleoecological data were dispersed across different resources, with differing
184 data architectures and degrees of accessibility, it has been difficult to synthesize data across
185 resources.

186 Neotoma builds on these prior efforts by providing 1) a consolidated and hence more
187 cost-efficient and sustainable IT structure, 2) an open and flexible data model based on
188 decades of experience with paleoecological data, and 3) a distributed and extensible
189 governance structure that promotes high-quality, curated data and establishes pathways for
190 new scientists and research groups to join and contribute. The Neotoma data model supports,
191 or can be extended to support, any kind of paleoecological or paleoenvironmental data from
192 sedimentary archives. Neotoma has focused on the Quaternary to Miocene section of the
193 geological record, and primarily supports research about ecological processes operating at
194 timescales of 10^2 to 10^6 years, but there is no hard limit to Neotoma's temporal extent. Neotoma
195 employs a distributed governance structure based upon Constituent Databases (see
196 *Governance and Data Use*), because the scientific expertise necessary for curating
197 paleoecological data is widely distributed across the scientific community. Each Constituent

198 Database corresponds to a particular dataset type (e.g. ostracodes, insects, vertebrates,
199 diatoms, pollen, organic biomarkers) or region and each is led by expert Data Stewards (list of
200 Stewards available at <http://bit.ly/2tzjEsZ>), with opportunities available for scientists who would
201 like to contribute their data and become Stewards, either for existing Constituent Databases or
202 to launch new ones. Neotoma also acts as a boundary organization and translator (Fig. F2,
203 Guston, 2001) among multiple interacting groups, bridging across the field- and lab-oriented
204 communities that contribute data to Neotoma, the diverse multiple research and educational
205 communities that use paleoecological data, and the informatics communities that build systems
206 for translating big data to knowledge.

207 Neotoma is part of the emerging field of paleoecoinformatics (Brewer et al., 2012), which
208 itself contributes to larger efforts in bioinformatics and geoinformatics to overcome bottlenecks
209 associated with data access, mobilize dark data, and maximize the power of scientific data
210 collected by networks of researchers (Ferguson et al., 2014; Hampton et al., 2013; Heidorn,
211 2008; Howe et al., 2008; Lynch, 2008). Related efforts include the Paleobiology Database
212 (<https://paleobiodb.org/#/>), the International Tree Ring Databank
213 (<https://www.ncdc.noaa.gov/data-access/paleoclimatology-data/datasets/tree-ring>), the
214 Paleoclimatology data holdings at the NOAA National Centers for Environmental Information
215 (<https://www.ncdc.noaa.gov/data-access/paleoclimatology-data>), Pangaea
216 (<https://www.pangaea.de/>), iDigBio and iDigPaleo (<https://www.idigpaleo.org/>), MorphoBank
217 (<https://morphobank.org/>), the Limnological Research Center and Continental Scientific Drilling
218 Office (<https://csdco.umn.edu/>), the Interdisciplinary Earth Data Alliance
219 (<http://www.iedadata.org/>), and The Strategic Environmental Archaeology Database
220 (<http://www.sead.se>). Interlinking efforts are also underway through e.g. NSF's EarthCube
221 program.

222 Neotoma also seeks to support and leverage on-going advances in paleosciences and
223 data sciences. In paleoecology and paleobiology, the rate of publications is increasing
224 exponentially (Uhen et al., 2013), which requires scalable informatic solutions; new
225 paleoecological proxies continue to emerge (e.g. organic biomarkers and compound-specific
226 stable isotopes, Bush and McInerney, 2013; Sachse et al., 2012; Zhang et al., 2006);
227 radiometric dating techniques and age modeling software continue to improve, enabling more
228 precise ecological inferences (e.g. Zazula et al., 2014); and the growth of data assimilation and
229 ecological forecasting approaches is requiring a closer and interactive coupling between
230 ecological data and mechanistic models (Dietze, 2017; Dietze et al., 2012). Relevant advances
231 in the data sciences include the development of efficient protocols (e.g. JSON, SPARQL) that

enable the establishment of networks of reliable and linked distributed resources across the internet (https://www.w3.org/blog/SW/2008/01/15/sparql_is_a_recommendation/) (ECMA International, 2013), the development of international standards for unique data identifiers (e.g. DOI, www.datacite.org), and the rapid advances in open-source and collaborative programming environments such as R, GitHub, or Jupyter that place state-of-the-art analytical tools at the fingertips of scientists (Goring et al., 2015; Ohri, 2014). Neotoma also supports the now-common requirement by funding agencies and journals for scientific data to be made publicly available (Nature, 2017).

Here we provide an overview of the Neotoma Paleoecology Database, its design principles, key concepts, software ecosystem, governance structure, and a snapshot of current data holdings and uploads. We summarize the structure of the Neotoma data model and the available capabilities for uploading, finding, exploring, analyzing, and validating data. A supplementary section points readers to additional information about the technical details of implementation, available in open-source code repositories (www.github.com/NeotomaDB), several on-line help manuals, and other publications (Goring et al., 2015; Grimm et al., 2014). We begin by describing key design principles and core semantic concepts that underlie the Neotoma data model, software design, and governance structure. Finally, we highlight ways that interested scientists can contribute to building this community resource, and we look ahead to current opportunities, challenges, and their solutions.

NEOTOMA DATABASE: DESIGN PRINCIPLES

Seven core principles and philosophies govern Neotoma's design.

1) Neotoma is a Spatiotemporal Paleoecological Database: The core mission of Neotoma is to store and openly share well-organized and curated information about the past occurrences and abundances of organisms, their geobiological signatures, and associated paleoenvironmental variables, in space and time (see *Data Holdings and Data Types* for a fuller listing). Informed interpretation of these data also requires Neotoma to store and curate information about the geographic and sedimentary characteristics of the field site and stratigraphic horizons from which fossils were collected or variables measured, age controls and age-depth models used to estimate time, and the identity and contact information of investigators. Other kinds of information are also relevant to informed paleoecological interpretation, but fall outside Neotoma's core mission and may be best curated by other communities and data resources, e.g. ecological traits, digital images of fossils, contemporary

264 genetic data, archaeological excavation data, reconstructions of past sea level and
265 paleogeography, and paleoclimatic simulations from earth system models. Neotoma's data and
266 governance models are designed to be flexible and extensible to other paleoecological or
267 paleoenvironmental proxies.

268 **2) Neotoma Consolidates IT, Distributes Scientific Governance:** Neotoma combines a
269 centralized database structure with a system of distributed scientific governance. All Neotoma
270 data are housed in a single relational database (see *Technical Specifications and Software*
271 *Ecosystem*) and organized according to a common set of core semantic concepts (see
272 *Neotoma Data Model: Fundamental Concepts*). The use of one data structure for multiple
273 paleoecological proxies reduces developer support costs, facilitates data discoverability and
274 reuse, and increases data interoperability. Neotoma comprises virtual Constituent Databases,
275 each encompassing a particular dataset type or region (e.g. North American Pollen Database,
276 European Pollen Database, North American Non-marine Ostracode Database, FAUNMAP,
277 Bugs) and each with its own Data Stewards (Fig. 3). Constituent Databases and allied
278 cyberinfrastructure resources can also create their own front-end portals into Neotoma through
279 the use of the Neotoma Application Programming Interface (API; see *Technical Specifications*
280 and *Software Ecosystem*). New Constituent Databases can be created to bring in new data
281 types and regions (see *Governance and Data Use Policy*) and new Data Stewards can be
282 readily trained (see *Governance and Next Steps* sections). This distributed governance
283 structure addresses the challenge of distributed scientific expertise and allows Neotoma to scale
284 as new records and proxy types are added to it.

285 **3) Neotoma is a Community-Curated Data Resource.** Neotoma's Data Stewards are
286 functionally similar to a journal's Board of Editors, charged with ensuring that data stored in
287 Neotoma conform to community-established data standards. Data input into Neotoma is led by
288 trained Data Stewards and Data Processors (see *Governance and Data Use Policy*) appointed
289 by their communities, with taxonomic names approved by Taxonomic Experts. Distributed
290 scientific governance is essential for Neotoma because no single individual or institution can be
291 expert in all the dataset types, regions, and time periods represented within Neotoma. Software
292 systems provide Stewards with automated tools to check for data inconsistencies, metadata
293 completeness, and taxonomic conformity with Neotoma standards prior to uploading (see
294 *Technical Specifications and Software Ecosystem*). This process of data validation,
295 standardization, and cleaning adds significant value (Lehnert and Hsu, 2015) and distinguishes
296 Neotoma from general-purpose and comprehensive data depositories such as DataDryad or
297 FigShare.

298 **4) Neotoma Data are Open:** Neotoma data are available to anyone with an internet
299 connection and are accessible through several interfaces (Fig. 4), each serving distinct user
300 communities. Neotoma uses a CC-BY license, allowing free reuse of data with proper
301 attribution to Neotoma and the original data contributors (see *Governance and Data Use Policy*,
302 <http://www.neotomadb.org/data/category/use>). By supporting open data, Neotoma prevents the
303 underuse and eventual loss of valuable paleodata that languish on individual computers
304 (Hampton et al., 2013; Heidorn, 2008) while promoting scientific transparency and
305 reproducibility across the community. Neotoma has an embargo policy for data contributed to
306 Neotoma prior to publication (*Governance and Data Use*) and is working on an embargo
307 management system, to promote good data management practice of organizing data early in a
308 project cycle, rather than at the moment of publication, when details may be difficult to recollect.
309 This embargo will also allow contributors to analyze their data in the context of the whole
310 database without releasing it for public access until publication. Neotoma supports data archival
311 and management plans required by many funding agencies and journals.

312 **5) Neotoma is a Living Database.** The life cycle of a paleoecological dataset does not
313 end with its first publication. Data are reused, as the original investigators and new teams
314 synthesize existing data to answer new questions. Errors may be caught and corrected during
315 subsequent syntheses (e.g. inaccurate transcriptions of geographic coordinates, inaccurate or
316 incomplete capture of all fossil data from a site, missing metadata). Derived inferences may be
317 updated with newer analytical methods (e.g. newer Bayesian age models, Blaauw and Christen,
318 2011; Parnell et al., 2008). Neotoma stores relatively stable raw data (e.g. number of fossil
319 specimens of a species in an assemblage, radiocarbon dates) as well as derived data (e.g. age
320 models), subject to change as scientific understanding advances. Of these, the most
321 changeable data tend to be age estimates and taxonomic names. Age estimates will change as
322 dating techniques and age-depth models improve. Taxonomic names may change if taxonomic
323 identifications of individual specimens are revised or if taxonomic nomenclature is revised or
324 updated.

325 **6) Community Engagement and Empowerment are Essential.** Neotoma serves diverse
326 communities (Figs. 2, 3). The constituent groups are essential to Neotoma's ability to grow and
327 scale upwards. Hence, Neotoma continually seeks to enlarge and support its community of
328 Data Stewards, contributors, and third-party developers (see *Governance and Next Steps*
329 sections). The sustainability of Neotoma is ultimately determined by the degree to which it
330 supports key research priorities of these users.

331 **7) Neotoma is Part of a Larger Ecosystem.** Many communities are gathering and
332 assembling their data, while others establish standards, vocabularies, and systems for sharing
333 data across systems. Data types are many and Neotoma's resources are few. Hence,
334 whenever possible, Neotoma will partner with other allied resources. For example, other
335 organizations have set standards for storing and representing information about individual
336 investigators (ORCID, <https://orcid.org/>, physical samples (IGSNs,
337 <http://www.geosamples.org/igsnabout>), geospatial data (Open Geospatial Consortium,
338 <http://www.opengeospatial.org/>), or funding agencies (FundRef registry built by CrossRef,
339 <https://www.crossref.org/services/funder-registry/>); Neotoma is adopting or moving towards
340 adopting these common standards. Similarly, Neotoma seeks to develop partnerships to
341 intersect paleoecological and paleoenvironmental data with other kinds of climatic,
342 archaeological, and ecological data (see *Next Steps*).

343 **NEOTOMA DATA MODEL: FUNDAMENTAL CONCEPTS**

344 This section describes high-level semantic concepts embedded within Neotoma's data model.
345 We first describe Neotoma's system for representing and storing the many kinds of sedimentary
346 sampling designs used by paleoecologists, then describe the kinds of information linked to
347 variables, and finally Neotoma's system for representing time. We do not attempt to describe in
348 detail Neotoma's relational database structure because most scientific users do not come into
349 direct contact with the actual relational database. However, many of the concepts described
350 here correspond to one or more data tables in Neotoma's relational database. For interested
351 scientists and developers, further information is available in the Neotoma Database Manual
352 (<http://www.neotomadb.org/uploads/NeotomaManual.pdf>).

353 **Sites, Collection Units, Analysis Units, Samples, and Datasets**

354 Paleoecological data from sedimentary archives have many commonalities: They
355 typically involve measurements of fossil organisms or *proxies* found in various geological
356 archives along some spatial direction, usually vertical *depth*, for which we estimate *time* with
357 uncertainty (see also Evans et al., 2013). These commonalities enable a common data model.
358 Within this general framework, many sampling systems exist, that vary within and among
359 subdisciplines and depositional environments. For example, paleolimnologists may collect one
360 or more sediment cores from a lake, with multiple kinds of measurements made on the cores
361 and subsamples from it; archaeologists may collect botanical or faunal specimens from surface

362 scatter or excavations; or vertebrate paleontologists may measure stable isotopes on bones
363 collected from a sediment section or cave deposit. In order to flexibly store data from these
364 different sampling methods, Neotoma uses a hierarchical arrangement of **sites**, **collection**
365 **units**, **analysis units**, **samples**, and **specimens**, with samples further grouped by type into
366 **datasets** (Fig. F5, Grimm et al., 2013).

367 A **site** is a geographic place from which paleoecological data have been collected.
368 Examples of sites include lakes, caves, archaeological excavations, and stratigraphic sections.
369 The spatial extent of sites is flexible, and tends to be defined based on field practice, ranging
370 from a single point to a lake or archaeological site with an extent measured in hectares. Key
371 properties include name, geographic coordinates, altitude, and areal extent. In Neotoma, the
372 spatial extent of sites is represented by bounding boxes with north and south latitudes and east
373 and west longitudes. The bounding box can circumscribe the site, (e.g. a lake) or may
374 circumscribe a larger area containing the site, either because site location is imprecisely known
375 (e.g. described as ‘on a gravel bar 5 miles east of town’) or because location is purposely kept
376 vague (e.g. to prevent looting and vandalism). Many legacy sites in Neotoma have point
377 coordinates. A site will have one to many collection units.

378 A **collection unit** is a place within a site from which a set of fossil specimens or samples
379 has been collected. Typical collection units include individual or composite cores from lakes and
380 peatlands, profiles from stratigraphic sections (e.g. river cutbanks, quarry walls), archaeological
381 or paleontological excavation features or contexts, isolated specimens (e.g. a bone collected
382 from a gravel bar), and surface samples. Collection units typically have spatial Cartesian or
383 geographic coordinates within a site. Collection unit properties include name, latitude-longitude
384 coordinates (represented as point coordinates with error), altitude, date of collection, and
385 metadata about collection methods and the depositional environment. The definition of
386 collection units and their spatial extent is flexible. For example, in a pit cave with three fossil-rich
387 sediment cones, each with several excavation squares, the collection units could be defined as
388 the individual squares, or as three composite collection units, one from each sediment cone. In
389 another example, consider a lake with two closely adjacent cores from the lake center, plus a
390 single core near the lake margin. In this case, the two central cores might be merged into one
391 composite collection unit, with a second collection unit representing the lake-margin core. A
392 collection unit will have one to many analysis units.

393 An **analysis unit** is a unique sampling location within a collection unit. Analysis units
394 are typically arrayed along a depth transect, usually oriented vertically. Analysis units may be
395 identified by depth and optionally thickness or by name and optionally ordinal position. Analysis

396 units may be arbitrary intervals or natural strata. Examples of analysis units include individual
397 depth intervals along a sediment core or within an excavation, individual strata or features within
398 an excavation, individual measurements taken along a transect from outer surface to inner core
399 of a speleothem, etc. Natural strata may be vertically superimposed, but may vary in depth and
400 thickness within a section or excavation, particularly in colluvial sections and sediment cones
401 below pit-cave openings. In this case, the “depths” may be ordinal positions or pseudo-depths.
402 However, in some cases, however, analysis units may be single-context features (e.g.
403 archaeological hearths and storage pits), which may be identified by name only, with no explicit
404 depth. An analysis unit will have one to many samples.

405 A **sample** is a single set of measurements of a single dataset type from an analysis unit.
406 Dataset types usually correspond to taxonomic groups that are the loci of scientific expertise
407 and Neotoma’s Constituent Databases, e.g. diatoms, ostracodes, pollen, plant macrofossils,
408 and terrestrial vertebrates. Samples and dataset types can also comprise geochemical
409 measurements (e.g. stable isotopes, organic biomarkers), physical measurements (e.g. loss on
410 ignition), or geochronological measurements (e.g. ^{14}C or ^{210}Pb dates). For example, the same
411 analysis unit from a stratigraphic section may have a vertebrate sample and a macrobotanical
412 sample; or an analysis unit from a sediment core may have pollen, diatom, ostracode, and
413 stable isotope samples. The analysis unit links together samples located in the same
414 stratigraphic interval, while datasets links samples from the same collection unit.

415 **Datasets** comprise all samples of the same data type from a collection unit. Datasets
416 are typically the subjects of publication and are a primary mode in which data within Neotoma
417 are packaged for delivery to users. For example, clicking on a site in Neotoma Explorer
418 (<http://apps.neotomadb.org/explorer/>) will return a list of all datasets at that site. Similarly, in the
419 *neotoma* package in R, Neotoma data are primarily passed to R in the form of datasets, using
420 the *get_datasets* and *get_download* functions (which respectively return dataset metadata and
421 data, Goring et al., 2015). Similarly, digital object identifiers (DOIs) are assigned to datasets,
422 but not to samples or analysis units (<http://data.neotomadb.org/datasets/5000/index.html>).

423 A **specimen** is the physical form of a biological object – e.g. a vertebrate bone or other
424 fossil – retrieved from a sample. Specimens are often curated and housed at museums,
425 geological surveys, or other repositories. Specimens often have catalogue numbers, accession
426 numbers, or other unique identifier assigned by their institution. In 2016-2017, the Neotoma
427 data model was extended to store specimen-level measurements, e.g. radiocarbon dates and
428 stable isotopic measurements from individual teeth or bones. Neotoma does not currently store
429 information about specimen morphometric traits, but the data model could be extended in this

430 direction, or Neotoma could link out to other databases that store specimen- or species-level
431 trait data (e.g. iDigBio, MorphoBank).

432 Variables

433 In Neotoma, variables store information about measured organisms or proxies of any
434 type. Variables have the property **taxon name** (equivalent to variable name) and the optional
435 properties **element**, **units**, and **context**.

436 In Neotoma, **taxon name** is used in the broad sense to include both organismal taxa
437 and physical “taxa” such as stable isotopes, organic biomarker compounds, and inorganic
438 minerals. Neotoma uses defined vocabularies of taxon names and, during the validation
439 process of uploading data to Neotoma, taxon names in uploaded files are automatically checked
440 and flagged if there is no match. New names may be proposed by Stewards and approved by
441 Taxonomic Experts (see *Governance and Data Use Policy*). Taxon names for organisms can
442 include non-Latin modifiers to indicate the level of uncertainty in the taxon identification. For
443 example, *Ambrosia*, *Ambrosia*-type, and cf. *Ambrosia* are three different taxa. The uncertainty
444 modifiers are included in the taxon name, rather than in a separate field, so as to indicate the
445 exact level of uncertainty and to faithfully record the original identification, in which the
446 uncertainty is usually included as part of the name. Thus, *Odocoileus* cf. *O. virginianus* indicates
447 that the genus identification is secure, but the species is uncertain; whereas, cf. *Odocoileus*
448 *virginianus* indicates that the genus identification is uncertain. This example might be the case
449 for regions in which *Odocoileus virginianus* (white-tailed deer) is the only *Odocoileus* species
450 biogeographically reasonable, thus if it is an *Odocoileus*, it must be *Odocoileus virginianus*. The
451 “cf.” designation may occur at any number of taxonomic ranks or cladistic nodes. With the
452 exception of the uncertainty modifiers (cf., aff., sp., spp., undiff., -type, ?), non-Latin modifiers
453 are included in parentheses, e.g. Poaceae (<50 µm), Leporidae (large), Eudicotyledoneae
454 (tricolpate, Hooghiemstra 1984 type 152). The forward slash symbol is used to indicate
455 identifications limited to a small number of taxa (usually two) with elements that cannot be
456 differentiated, e.g. *Ostrya/Carpinus*. This packaging of taxon name and uncertainty into a single
457 field carries the advantage of staying true to paleontological tradition and nomenclature, but
458 hinders integration with contemporary biodiversity databases.

459 If taxonomic names used by data contributors or publications are changed, the original
460 name can be stored. Database-wide changes to taxon are always stored, for example a name
461 change due to taxonomic revision. Different constituent databases have somewhat differing
462 practices for nomenclatural synonymizations made at the time of initial data validation and

463 upload: some retain the original name, others replace with the currently accepted synonym. For
464 homotypic synonymizations that do not involve a change in circumscription (e.g. Gramineae ≡
465 Poaceae), original names do not have to be stored; but for heterotypic synonymizations or any
466 name change that involves a potential change in circumscription, original names should be
467 stored. In addition, name changes due to re-identification should be stored, for example a
468 specimen initially identified as *Mammuthus* but later re-identified as *Mammut*.

469 For organismal taxa, the **element** is the organ or part of the organism that was identified.
470 Thus, for pollen datasets, element names include 'pollen,' 'spore,' and 'stomate.' For plant
471 macrofossil datasets elements include 'leaf,' 'seed,' 'bud scale,' 'microstrobilus,' and 'wood.' For
472 fossil insects, elements include 'heads,' 'pronota' (thoraces), 'elytra' (left and right), and
473 'aedeagii' (reproductive organs). For vertebrates, the element is generally the bone or tooth
474 identified, such as 'femur' or 'tooth, third molar.' Elements, particularly vertebrate elements, may
475 have components, including symmetry, portion, and maturity, entered in that order and
476 separated by semicolons, e.g. 'femur;left;distal;fused,' or 'tooth, third molar;lower left.' Elements
477 are also used for modifying physical variables. For example, for the variable loss-on-ignition, the
478 element is the temperature, e.g. 500°C. Neotoma also uses defined vocabularies for elements.

479 **Units** are the measurement units in which the variable is measured. For organismal
480 taxa, the most common units are NISP (Number of Identified Specimens, often called a "count"
481 for microfossils), MNI (Minimum Number of Individuals), and presence. For example, five left
482 femurs would indicate at least five organisms (MNI=5), but five pollen grains could have one to
483 five source plants (NISP=5). For presence data, a value of 1 is entered for presence or the cell
484 is left blank. Zeroes or absences are not stored in Neotoma for organismal taxon variables,
485 because the true absence of a taxon is difficult to definitively establish, given that probability of
486 detection is a function of sampling effort, ecological rarity, and taphonomic processes (Birks and
487 Line, 1992; Olszewski and Kidwell, 2007; Weng et al., 2006). However, partial evidence for
488 taxonomic absence exists for some samples if that taxon was identified in other samples in the
489 same dataset, under the assumption that the analyst looked for that taxon in all samples in a
490 dataset. Neotoma also allows semi-quantitative systems for measuring abundance, e.g. the 1-5
491 relative abundance scale often used in the rodent midden literature (Spaulding et al., 1990) or
492 other relative scales used in archaeobotany. Geochemical and physical variables can have
493 many kinds of measurement units, e.g. 'percent,' 'per mille,' 'meq/L,' 'mg/L,' and so on. For
494 physical and geochemical measurements, such as $\delta^{13}\text{N}$, for which zero can be a measured
495 value, zeroes are stored. Neotoma uses defined vocabularies for units.

496 **Context** refers to a depositional context that may influence the interpretation of the
497 taxon. Examples include anachronic, redeposited, or intrusive, which imply that the taxon was
498 deposited at different time than its sediment matrix. A Cretaceous pollen grain may be reworked
499 and redeposited into more recent sediments; a modern *Sus* (pig) bone may be intrusive in
500 Pleistocene sediment. Anachronic simply implies the taxon is of a different age. Contexts
501 sometimes used with pollen are clump and anther, where clumps of pollen or anther fragments
502 may indicate an over-representation of the taxon, e.g. a bee carrying usually infrequent
503 entomophilous pollen may have fallen into the sediment. It is possible for fossils in a dataset or
504 sample to be from the same taxon but have different contexts, e.g. both clumped and non-
505 clumped pollen of one taxon, or Holocene and older reworked Betulaceae pollen
506 (distinguishable by preservational differences) in the same assemblage.

507 **Time, Age Controls, Relative Ages, Age-Depth Models, Chronologies**

508 In Neotoma, the age estimates attached to samples and specimens are treated and
509 stored as a derived variable, that must be estimated through a combination of absolute age
510 controls, relative age controls, event stratigraphic ages, and age-depth models fitted to those
511 controls (e.g. Parnell et al., 2008). Some types of geochronological information are more stable
512 than others. For example, individual geochronological measurements will remain stable, but the
513 set of age controls available at a site or collection unit can change as new dates are obtained.
514 The fitted age-depth models and derived age inferences are also changeable, as estimates of
515 radiometric decay constants are updated, radiocarbon calibration curves are adjusted, the
516 quality of individual age controls is reassessed, new statistical age-depth models are developed,
517 etc. Hence, the Neotoma data model separately represents and stores these layers of
518 information about time as *age controls*, *relative ages*, *age-depth models*, and *chronologies*.
519 Most of the definitions described here are originally from Grimm et al. (2014).

520 An **age control** is an estimate of absolute age, often with a specified uncertainty, for a
521 sample within a core or stratigraphic profile that is used to constrain an age model for that core
522 or profile. Examples of age controls include radiocarbon and other radiometric dates, optically
523 stimulated luminescence, biostratigraphic events, tephras, core top, coins or other dated cultural
524 artifacts, etc. (Blois et al., 2011; Giesecke et al., 2014). Age controls are primary data in
525 Neotoma and generally assumed to be fixed and unchanging; barring data entry error: new age
526 controls can be added, but existing age controls are not modified. Age control data are stored
527 in tables that are separate from but linked to tables that store information about age-depth

528 models and chronologies (see below). In Neotoma, age control data are stored in the
529 Geochronology table.

530 Radiometric ages are the most common kind of age control stored in Neotoma, and
531 radiocarbon dates are the most common kind of radiometric date (17,054 of 18,483 age controls
532 in Neotoma are radiocarbon dates, as of November 7, 2017). Radiocarbon dates are stored in
533 original radiocarbon years, with counting uncertainties stored as one standard deviation.
534 Calibrated ages are stored in Neotoma as components of chronologies (below), or simply
535 regenerated by users as needed. Other radiometric metadata include depth and thickness of
536 sample, material dated, lab identifier, $\delta^{13}\text{C}$ (for radiocarbon dates), instrumental measurement
537 system, and publication information.

538 **Relative ages** store information about the association of analysis units with formally
539 recognized relative age scales based on the stratigraphic record or a series of geological
540 events. These formations and events have their own age estimates, which can change over
541 time. A relative age encompasses a range of time (i.e. it has an upper and lower age bound),
542 and samples assigned to a relative age event are assigned those ages. Examples include
543 Marine Isotope Stages, Heinrich stadials, geomagnetic chronos, archaeological periods, and
544 North American land mammal ages. Thus, a sample assigned to Marine Isotope Stage 5e would
545 be assigned a sample age of 130-116 ka based on the current authoritative estimate (currently,
546 Lisiecki and Raymo, 2005). Optionally, this sample age could be further constrained by other
547 criteria. A full list of relative age scales is available via the Neotoma API
548 (<http://api.neotomadb.org/v1/dbtables/RelativeAgeScales>, in JSON format). Currently these
549 relative age scales and associated age estimates are stored in Neotoma in the RelativeAges
550 table and must be updated by Stewards; ideally these ages would be dynamically updated by
551 linking to other authoritative data resources on stratigraphic age.

552 **Event stratigraphic ages** are globally synchronous, single-event stratigraphic markers,
553 ideally with an age and error, which can be used in age models similar to geochronological
554 ages. Volcanic deposits tephras often serve this purpose. For example, the Mazama tephra,
555 which occurs over a large area of western North America, has been dated to $7,627 \pm 150$ cal yr
556 BP (Zdanowicz et al., 1999). The Hekla 1104 tephra originating from Iceland, but distributed as
557 far east as Ireland, is historically documented to AD 1104 or 846 cal yr BP (Boyle, 1999).
558 Some event stratigraphic ages are the boundaries of relative age units, e.g. geomagnetic
559 polarity reversals are event stratigraphic ages and are also the boundary ages for geomagnetic
560 chronos, which are relative ages. Hence, a fossil vertebrate assemblage might occur within a unit
561 stratigraphically assigned to the reversed polarity Matuyama chron, with a relative age of 2.581-

562 0.781 Ma, while the Brunhes/Matuyama geomagnetic polarity reversal, which might be an age
563 control for an age model in a core or section, is dated to 0.781 Ma (Ogg and Smith, 2005).
564 Event stratigraphic ages are treated in the Neotoma data model as instantaneous; although, of
565 course, the actual events occurred over periods of time, from weeks to a few years for a
566 volcanic eruption to perhaps several hundred or even a few thousand years for a geomagnetic
567 reversal (e.g. Clement, 2004). However, for most practical applications, this error should be
568 relatively small and ideally incorporated within the age error estimate.

569 An **age-depth model** is an algorithm used to estimate the age-depth relationship for a
570 given stratigraphic profile based on the age controls available for that profile and prior
571 knowledge. Age-depth models are used to estimate ages for depths not directly associated with
572 an age control or to resolve discrepancies among age controls. Examples of age-depth
573 modeling programs include classical age modeling approaches (linear interpolation, linear
574 regression, polynomials, and splines, Blaauw, 2010) and Bayesian approaches such as Bacon
575 (Blaauw and Christen, 2011) or BChron (Parnell et al., 2008). In Neotoma, information about
576 age-depth models is stored in the Chronologies table.

577 A **chronology** is a series of estimated ages for a set of samples in a collection unit,
578 ideally with associated uncertainty estimates. Chronologies usually derive from an age-depth
579 model and a set of age controls (radiometric, relative, events, or other) used to constrain that
580 model. In Neotoma, information associated with chronologies is parsed into three tables: 1) The
581 ChronControls table, which stores information about the age controls used to constrain an age-
582 depth model. These age controls can, but do not have, to correspond to the age controls stored
583 in the Geochronology and other primary data tables; this flexibility allows scientists to remove
584 age controls deemed to be inaccurate or add other age constraints to the age-depth model (e.g.
585 assigning a modern depositional age to the top of the stratigraphic profile). 2) The Chronologies
586 table stores information about the age-depth model, its parameters, the analyst, and other
587 metadata. 3) The SampleAges table stores the resultant inferred ages for individual samples.
588 All Neotoma chronologies have a unique identifier and are also linked to a specific collection
589 unit.

590 Chronologies can be stored as calendar years before present, radiocarbon years before
591 present, or calibrated radiocarbon years before present. The chronology and associated age-
592 depth model originally published for a collection unit are stored provided that the age controls,
593 sample ages, and metadata sufficient to replicate the age model are published or provided by
594 the contributor. A collection unit may have multiple chronologies or none, if no age information
595 is available for the collection unit. Hence, each sample in the collection unit may have multiple

596 age estimates, each linked to a unique chronology. New chronologies may be added to
597 Neotoma, and we envision this component of the Neotoma age model to be dynamic over time,
598 as users download data and build new chronologies.

599 To handle the multiplicity of chronologies, Neotoma allows one chronology per collection
600 unit to be designated as a default chronology. Choice of default chronologies is made by Data
601 Stewards, and this choice can be revised as new chronologies are added. Default chronologies
602 may be stored in calendar, radiocarbon, or calibrated radiocarbon years before present.
603 Chronologies in radiocarbon years are not recommended, but exist as legacies in Neotoma.
604 We anticipate that over time, these default radiocarbon chronologies will be replaced by
605 updated chronologies in calibrated radiocarbon years. Researchers using Neotoma data are
606 cautioned to critically examine its chronologies and encouraged to contribute new and revised
607 chronologies.

608 **DATA HOLDINGS AND DATA TYPES: CURRENT STATUS AND TRENDS**

609 As of November 8, 2017, Neotoma holds over 3.8 million data records from 17,275
610 datasets and 9,269 sites (Fig. 6). Each data record is the measured value of a single taxon or
611 other variable from a single sample. The Neotoma taxa table is a dictionary with over 29,000
612 taxa, fossil morphotypes, geochemical variables, and other variable names.

613 Data volumes in Neotoma have been rapidly growing as data are uploaded, with a 30%
614 increase since 2014, when a new wave of data uploads began, following the extension of the
615 Tilia software package (see *Technical Specifications and Software Ecosystem* below) was
616 extended to enable data validation and upload to Neotoma (Fig. 6). As of July 13, 2017,
617 Neotoma holds 2,954 pollen datasets (3,274,501 data records), 2,600 pollen surface sample
618 datasets (56,205 data records) 3,669 vertebrate fauna datasets (59,278 data records), 388
619 diatom datasets (238,344 data records), 637 diatom surface sample datasets (29,968 data
620 records), 554 ostracode surface sample datasets (2,410 data records), 384 macroinvertebrate
621 datasets (805 data records), 283 plant macrofossil datasets (10,654 data records), 177 insect
622 datasets (19,766 data records), and a number of other dataset types (Table 1). Of these, 5,226
623 datasets are from sediment cores, 562 from rodent middens, 285 from excavations, and 1,051
624 from stratigraphic sections. Neotoma holds 3,842 geochronologic datasets, with 18,543
625 individual geochronologic measurements. These data volumes have made Neotoma one of the
626 largest structured repositories of geochronological data.

627 Although Neotoma has focused primarily on paleoecological and geochronological data,
628 it also stores associated physical and geochemical proxies. It includes modern water chemistry
629 data (1,317 datasets, 18,487 data records), loss-on-ignition data (190 datasets, 24,903 data
630 records), charcoal (100 datasets, 14,428 data records), and physical sediment measurements
631 (55 datasets, 1292 data records). Trial datasets have been uploaded for X-Ray Fluorescence
632 (XRF) and X-Ray Diffraction (XRD), and sedimentary geochemistry data. Surface samples can
633 be flagged in Neotoma during upload, to facilitate their use in the building of modern calibration
634 datasets for transfer functions (Birks, 1995). Some associated measurements of environmental
635 variables (relating to water chemistry) are enabled for the ostracode and diatom surface
636 samples stored in Neotoma, and we plan to extend the Neotoma data model to store other
637 environmental variables (e.g. climate variables) associated with surface samples. The Neotoma
638 data model was recently extended to include stable isotopic data and metadata for $\delta^{18}\text{O}$, $\delta^{13}\text{C}$,
639 $\delta^{15}\text{N}$, δD , $\delta^{34}\text{S}$, and $^{87}\text{Sr}/^{86}\text{Sr}$. Isotopic measurements can be stored for samples from
640 sedimentary profiles or from individual fossil specimens. We are expanding the Neotoma data
641 model to store organic biomarker data and taphonomic measurements on vertebrate fossils.
642 Because the formation of Constituent Databases in Neotoma is a voluntary, bottom-up process
643 (*Governance and Data Use Policy*), whether to extend Neotoma to other proxies largely
644 depends on 1) interest by paleodata communities in using Neotoma to house their data and 2)
645 developer and Steward time to extend data tables and metadata variables as needed.

646 Note that Neotoma currently focuses on storing primary paleoecological measurements
647 and generally does not emphasize the storing of derived inferences. For example, Neotoma
648 currently does not store indices of community diversity (richness, evenness, etc.), paleoclimatic
649 reconstructions, biomes or other paleovegetation reconstructions, etc. Our general philosophy
650 is that, given finite resources and the rapid pace of generating these inferences by the scientific
651 community, these reconstructions are best generated and managed outside of Neotoma, ideally
652 using workflow methods that clearly link all derived inferences back to the primary data
653 resources inside Neotoma. However, this boundary between primary data and secondary
654 inferences is not set in stone, and could be revised depending on research and development
655 priorities within and among Constituent Databases and their research communities.

656 **TECHNICAL SPECIFICATIONS AND SOFTWARE ECOSYSTEM**

657 The Neotoma software ecosystem (Fig. 4) has Neotoma's relational database at its heart and
658 includes multiple systems for finding, exploring, visualizing, downloading, processing, and
659 uploading data. We describe each in turn.

660 **Database**

661 The Neotoma Database is currently hosted on servers at Pennsylvania State University,
662 maintained by the Center for Environmental Informatics (CEI). Neotoma uses a relational
663 database structure that was originally deployed in Microsoft SQL Server and now is being
664 migrated to PostgreSQL, so that it will rely primarily on open source technology. Neotoma's
665 relational database structure continues to evolve over time, as new data types and metadata
666 fields are added. A description of the relational database structure and tables is available in an
667 on-line manual (<http://neotoma-manual.readthedocs.org/en/latest/>).

668 Data uploaded to the Neotoma relational database are protected by multiple backup
669 measures at CEI, including redundant disk storage, off-site mirroring, file system snapshotting,
670 regular tape backup, and duplication of the backup set. Complete snapshots of the Neotoma
671 database are posted to the Neotoma website (<https://www.neotomadb.org/snapshots>) and to
672 figShare (e.g. <https://dx.doi.org/10.6084/m9.figshare.3376393.v1>). figShare in turn ensures long
673 term data sustainability through its partnership with the Digital Preservation Network who obtain
674 periodic snapshots of the entire figShare collection (including the Neotoma snapshots), and
675 replicate it across at least two Replicating Nodes, including the Academic Preservation,
676 DuraCloud Vault, Stanford Digital Repository, Texas Preservation Node, and HathiTrust.
677 Neotoma also will send these snapshots to the Paleoclimatology branch of the NOAA National
678 Center for Environmental Informatics. Database snapshots are intermittent at present, but we
679 plan to establish a regular automated schedule on a quarterly frequency. Neotoma is a certified
680 member of the International Council for Science's World Data Service (ICSU-WDS), which sets
681 standards for open, quality-assured, and sustained stewardship of scientific data.

682 **Data Retrieval: Finding, Exploring, and Downloading**

683 Multiple avenues exist to find, explore, and obtain data from Neotoma, each serving
684 different needs and users. Neotoma data can be accessed via the Neotoma Explorer
685 interactive map-based interface (<https://apps.neotomadb.org/Explorer/>), through a RESTful API
686 (<https://api.neotomadb.org/>), and through digital object identifiers (DOIs) that provide persistent

687 and unique identifiers for every dataset in Neotoma (e.g.
688 <https://data.neotomadb.org/datasets/1001/>). Further, Neotoma data can be searched, viewed,
689 and analyzed through Neotoma Explorer and the Neotoma API, and also the *neotoma* R
690 package and stratigraphic and map-based visualizations in Tilia. Through partnership with the
691 Earth Life Consortium (<http://earthlifeconsortium.org/>), we are developing wrapper APIs that can
692 simultaneously search for paleobiological data in Neotoma, the Paleobiology Database, and
693 other partner databases (http://www.earthlifeconsortium.org/api_v1/ui/). A full set of links,
694 manuals, code repositories, and other resources are provided in the Supplementary Material.

695 **Neotoma Explorer** (<https://apps.neotomadb.org/explorer/>), a map-based web
696 application for searching, visualizing, and downloading data, is Neotoma's primary data
697 discovery portal. Users can generate flexible queries for properties such as taxon name,
698 variable type, time window, location, constituent database, site name, and researcher name.
699 Results are displayed on an interactive map, and users can quickly retrieve site and dataset
700 metadata by clicking on sites of interest. Users interested in deeper exploration of datasets can
701 then view them in Explorer's Stratigraphic Diagrammer to inspect age models, stratigraphic
702 plots, and associated publications. Datasets can be downloaded as delimited text files or shared
703 via links that make use of Neotoma's API (e.g.,
704 <http://apps.neotomadb.org/Explorer/?datasetid=1768>). Searches can be saved as JSON files,
705 which can be archived or shared with other users, who can reopen them in Neotoma Explorer
706 (by dragging the JSON file to the Explorer window) to redisplay the results and map
707 configuration produced from the query. Following our design principle of promoting openness,
708 Explorer is built on the Open Web Platform using HTML5, CSS3, and JavaScript, enabling
709 cross-browser and cross-platform support, and makes use of the open-source Dojo and
710 OpenLayers libraries (Roth et al., 2014).

711 **APIs** enable programmatic access to the database by third-party developers and
712 software applications. Likely API users include: scientists who need to incorporate the most
713 current data into analytical workflows such as scripts written in R; organizations that want to
714 distribute Neotoma data (with attribution) via their own data portal or web interface (e.g.
715 customized and branded web portals for individual Constituent Databases); and developers
716 creating standalone applications for data analysis and display. Known third-party users of
717 Neotoma APIs include NOAA Paleoclimatology (<https://www.ncdc.noaa.gov/data-access/paleoclimatology-data>), which has a data search portal that will search and retrieve data
718 housed in either NOAA or Neotoma; Flyover Country (<http://fc.umn.edu/>), a mobile app-based
719 program for travelers to discover geological data and knowledge during their journeys (Loeffler
720

721 et al., 2015); and the Global Pollen Project (<https://globalpollenproject.org/>) (Martin and Harvey,
722 2017), a community platform for pollen identification. The APIs are implemented as platform-
723 and language-independent RESTful web services; response formats include JSON and XML.

724 The R package ***neotoma*** uses the Neotoma API to pass data into R for further analysis
725 (Goring et al., 2015). R is an open-source fourth-generation programming language for
726 statistical analysis and graphics. Many paleoecological statistical and visualization packages
727 have been developed for R, such as *analogue* (Simpson, 2007; Simpson and Oksanen, 2015),
728 *rioja* (Juggins, 2015), *bchron* (Parnell et al., 2008), *bclim* (Parnell et al., 2016), *clam* (Blaauw,
729 2010), and *bacon* (Blaauw and Christen, 2011). Development of *neotoma* is ongoing (v1.7.0 is
730 available via the Comprehensive R Archive Network, <https://cran.r-project.org/web/packages/neotoma/index.html>) and the living *neotoma* code, feature requests,
731 bug reports, and development are open and available through GitHub
732 (<https://github.com/ropensci/neotoma>).

733 **DOIs.** We have recently begun assigning persistent and unique digital object identifiers
734 (DOIs) to all Neotoma datasets, which will facilitate citation of Neotoma data and linked-data
735 systems for sharing and connecting earth science data (Duerr et al., 2011). DOIs are created
736 using the EZID system (<https://ezid.lib.psu.edu/>), affiliated with DataCite
737 (<https://www.datacite.org/>), through a license with the University of Wisconsin-Madison Library
738 and University of Illinois. DOIs are assigned at the level of datasets, and each DOI has a
739 landing page (e.g. <https://data.neotomadb.org/datasets/1001/>) that is designed to provide
740 information to both human and machine users. The DOI pages point to living versions of the
741 Neotoma data; if the Neotoma data are updated, these changes will be automatically detected
742 and incorporated in the DOI landing pages. This approach differs from e.g., Pangaea
743 (www.pangaea.de), in which DOIs point to static versions of datasets.
744

745 **Data Validation, Upload, and Management**

746 All data added to Neotoma are reviewed by a Data Steward (see *Governance and Data*
747 *Use*) before upload; this expert curation is central to Neotoma's mission of providing high-quality
748 scientific data. Data entry, curation, and upload into Neotoma are handled through the ***Tilia***
749 software (<https://www.tiliait.com/>). Tilia was originally developed as a DOS program to visualize
750 and analyze pollen stratigraphic data (Grimm, 1988) and is still often used for this purpose. Tilia,
751 now a Windows program, still maintains its end-user functions for managing, analyzing, and
752 visualizing stratigraphic data, but its capabilities have been extended to support validation and
753 upload of data to Neotoma, as well as direct download of data from Neotoma. Only Neotoma

754 Stewards who have password access to Neotoma have access to the upload capability; any
755 user can download data. Tilia now generates .tlx files in extensible markup language (XML)
756 format, which is a simple, extensible text-based markup language format. Tilia contains multiple
757 validation procedures for quality control during upload, including checks against controlled
758 vocabularies, for missing or duplicate data, for data inconsistencies, common errors, and
759 commonly omitted metadata.

760 Neotoma uses controlled vocabularies for taxa, elements, units, contexts, geopolitical
761 units, depositional environments, geochronological measurements, chronological controls, and
762 other variables that store text-based names. These names are stored in Neotoma tables, and
763 the Tilia validation process checks names in .tlx files against these names. Central to the
764 database is the Taxa table, which contains the names of all taxa in the database. Names not
765 found in the Taxa table during Tilia validation are listed, and Stewards can check for spelling
766 errors or formatting differences for non-Latin modifiers by searching with wildcards. If the name
767 is valid but not yet in Neotoma, the Steward can upload it to the database. Taxonomic names
768 are separately validated by designated taxonomic experts. If the Steward is a taxonomic expert,
769 then the date of entry is recorded as the date of validation. Otherwise the field for validation date
770 is left empty until the taxon is validated. Tilia has a tool that shows the Neotoma taxa in a
771 taxonomic hierarchy, easily allowing the Steward to place the new taxon in its correct taxonomic
772 position (Fig. 7).

773 Other Tilia validation steps include checking for valid dataset types, for valid
774 combinations of taxon and element (e.g. disallowing the combination of 'Picea' and 'femur'), that
775 latitude-longitude coordinates for collection units fall within site bounding boxes, that elevations
776 fall within possible limits, and that the younger/older reliable age bounds for chronologies are
777 not reversed (a common error). To validate longitudes and latitudes, the Steward is shown a
778 world map with the hemisphere indicated, which must be approved. Positive $\delta^{13}\text{C}$ values
779 associated with radiocarbon dates (a common entry error) are flagged. Tilia will ask the Steward
780 whether the top sample of a stratigraphic sequence should be flagged as a modern surface
781 sample. During the validation process, error messages, warnings, or notes may be issued.
782 Notes indicate omissions of optional metadata items that are not required and often non-existent
783 or unavailable. Warnings indicate omissions of optional metadata items that nevertheless are
784 highly desirable. Warnings are also issued for likely or possible inconsistencies, but which
785 nevertheless may be correct. Errors must be addressed before upload is possible. Once data
786 are validated with no errors, authorized Data Stewards can upload data from Tilia directly to
787 Neotoma through password-protected web services.

788 Tilia has a customized API for data upload and download. Controlled vocabularies that
789 appear in Tilia dropdown pick lists, such as taxa and geopolitical names, are held in local XML
790 lookup files, which should be synchronized periodically from their counterparts in the central
791 Neotoma database. Datasets downloaded from Neotoma to Tilia include all relevant metadata,
792 including various notes and comments that may have been entered. This functionality enables
793 enhanced review, visualization, and analysis of Neotoma data beyond that possible in Neotoma
794 Explorer. The downloaded data also facilitate training for end users and Data Stewards, by
795 providing model datasets in Tilia format. The Data Steward version of Tilia also allows Data
796 Stewards to amend data already in Tilia. For example, metadata items such as latitude-
797 longitude coordinates may be corrected, or missing metadata items for sites and collection units
798 may be added. Publications can be corrected or added. Contact information can be updated.
799 Datasets can be added to sites already in Neotoma. New chronologies and sample ages can be
800 added. The ultimate goal is to enable Data Stewards to correct or add to any data or metadata
801 item within their Constituent Database through the Tilia interface. Changes made through Tilia
802 are logged, Neotoma preserves snapshots of prior database versions, and we are collaborating
803 with others on building an annotation system, with support from NSF's EarthCube program. The
804 ability to upload and amend data is a significant power available to Data Stewards, hence the
805 need for experts to serve as Data Stewards and for Data Stewards to act judiciously when
806 modifying data and metadata.

807 Tilia software is available from <https://www.tiliait.com/>. The free version has all the
808 spreadsheet and metadata form options, including the ability to download data from Neotoma,
809 which can be copied to other spreadsheet programs. A licensed version for Stewards is
810 available at no cost and exposes additional options for data visualization, validation, upload, and
811 management, with a password needed for any action that may alter the database. Other
812 licensed versions of Tilia with graphics capabilities are available and are priced to cover
813 software licensing costs associated with Tilia development.

814 The Tilia workflow for uploading data follows a model in which site-level datasets and
815 associated metadata are uploaded individually. For larger data ports, where many site-level
816 datasets must be exported from one database into Neotoma (e.g. from the Access tables
817 storing data in FAUNMAP and the European Pollen Database), we have taken two approaches.
818 One is to place the database on the Neotoma server, and write customized SQL procedures
819 and web services to download datasets directly to individual Tilia files for validation and upload
820 to Neotoma. The second is to write customized scripts in R or Python that export datasets from
821 the other database to individual Tilia .tlx files, which can then be opened for validation and

822 upload by Data Stewards. Examples of these batch export scripts, developed for FAUNMAP
823 are available on GitHub (https://github.com/NeotomaDB/FAUNMAP_Import).

824 GOVERNANCE AND DATA USE

825 Governance and data use policies are designed to support Neotoma's core goals of data
826 openness and distributed scientific governance. In particular, these policies are intended to: 1)
827 Make Neotoma data open and available to all interested scientific and public communities; 2)
828 Build a governance structure that accommodates both a centralized cyberinfrastructure and a
829 highly distributed scientific community of expertise; and 3) Empower and facilitate individual
830 Data Stewards and Constituent Databases to set data acquisition priorities, curate data, and
831 establish data quality standards and nomenclatures.

832 All Neotoma data are free to use through a CC BY 4.0 license. Complete attribution of
833 Neotoma data includes a reference to the Neotoma Paleoecology database, constituent
834 databases where relevant, and references to all original investigators and publications. An
835 embargo policy has been developed and included in the Neotoma website and technical
836 implementation is underway. See the Neotoma website
<https://www.neotomadb.org/data/category/use>) for a full description of the Neotoma data use
838 policy and for the data use statements and citation formats for specific Constituent Databases.

839 With respect to governance, key needs include 1) an extendible, scalable governance
840 structure that is easily open to new members, 2) effective executive decision-making and
841 responsibility that is bounded by community oversight, and 3) mechanisms to ensure that
842 Neotoma is curated by a community of professionals and scientific experts. Here we briefly
843 summarize the main elements of Neotoma's governance (Fig. 3) and data use policies. A full
844 description of Neotoma governance is described in its bylaws, available at
<http://www.neotomadb.org/about/category/governance>.

846 Neotoma is governed by a **Leadership Council** (NLC) that sets policy and represents
847 the scientific perspectives of Constituent Databases and their Data Stewards (Fig 3).
848 Councilors serve for four-year renewable terms, with one-fourth of the NLC up for election each
849 year. The Council is elected by Neotoma's **members**, who are professional researchers and
850 educators who contribute to and use Neotoma data and are interested in helping govern the
851 database. Membership can be requested by any individual through a simple webform
852 (<https://tinyurl.com/NeotomaMember>), and requests are approved by Neotoma's Nominations

853 and Membership Working Group, chaired by Neotoma's Associate Chair. Data Stewards are
854 automatically granted Neotoma membership.

855 The NLC delegates responsibility for day-to-day operations to an **Executive Team** (Fig.
856 3), consisting of an Executive Chair, Associate Chair, and two other members. All positions on
857 the Executive Team serve staggered four-year terms and are selected from and by the NLC.
858 Other teams within the NLC include the Education and Outreach Working Group, the
859 Informatics and Technology Working Group, and the International Partnerships Working Group.

860 **Constituent Databases** are a core concept in Neotoma, and a mechanism by which
861 scientific data governance is distributed among the multiple fields of scientific expertise that
862 Neotoma data embody. (A list of Neotoma's Constituent Databases can be obtained through
863 the Tilia API - <https://tilia.neotomadb.org/retrieve/?method=GetConstituentDatabases> - which
864 returns a JSON object). All data in Neotoma are associated with a Constituent Database, each
865 curated by a community of Data Stewards, Taxonomic Experts, and Data Uploaders.
866 Constituent Databases are responsible for vetting and uploading their community's data to
867 Neotoma, setting priorities and quality standards for data uploads, managing taxonomic names,
868 and appointing and training their Data Stewards and Taxonomic Experts. Constituent
869 Databases may develop variants on the standard Neotoma Data Use policy (see below). Some
870 Constituent Databases may be active for a few years, during e.g. a data mobilization campaign
871 linked to specific large-scale synthesis project and research grant. Other Constituent
872 Databases may be active indefinitely, when individuals and communities use Neotoma as their
873 platform for archiving, managing, and sharing data.

874 **Data Stewards** serve a role analogous to that of editors in a peer-reviewed scientific
875 journal. Much of the day-to-day power and responsibility devolves to Stewards for ensuring
876 high-quality data uploads to Neotoma. Data Stewards are authorized to upload data directly to
877 Neotoma within their Constituent Database, and to modify data within their Constituent
878 Database. Data Stewards often work with one or more **Data Processors**, often students or
879 other assistants who assist in the preparation and entry of data into Tilia for eventual upload to
880 Neotoma. **Taxonomic Experts** are a type of Data Steward that can also authorize the addition
881 of names to Neotoma's list of accepted taxa names and variables.

882 **NEXT STEPS**

883 **New Users, Contributors, and Communities**

884 Neotoma welcomes new members, data contributors, users, data stewards, and
885 constituent databases. There are many avenues for participation by interested scientists. One
886 simple step is to become a Neotoma member (<https://tinyurl.com/NeotomaMember>). Scientists
887 interested in contributing their data to Neotoma should contact a Data Steward
888 (<https://www.neotomadb.org/data/category/contribution>). Research labs with a lot of data may
889 want to consider steward training for someone in their lab. Training webinars can be scheduled
890 by request and are led by current Stewards. For scientists interested in learning how to access
891 and use Neotoma data, we have posted learning materials on-line and periodically hold user-
892 oriented training workshops (see Supplementary Materials).

893 Research teams interested in building regional- to global-scale data syntheses, both
894 paleoecological and paleoclimatic, may find Neotoma useful as a data synthesis platform. In
895 some cases, these efforts could lead to targeted data mobilization campaigns, uploads of data
896 to Neotoma, and the chartering of new Constituent Databases.

897 For example, the ANTIGUA project (e.g. Barnosky et al., 2016) is using Neotoma to
898 store fossil occurrences and age constraints for South American megafaunal species, with
899 datasets currently being processed for upload to Neotoma. The SKOPE project, studying
900 human-environment interactions in the southwestern US (e.g. Bocinsky and Kohler, 2014), is
901 using Neotoma as a platform for accessing paleoecological data and, in the process,
902 discovering data corrections and additional records for addition to Neotoma. PalEON,
903 interested in understanding climate-driven vegetation dynamics over the last 2,000 years, has
904 been discovering new records for addition to Neotoma and updating age models as part of its
905 development of the STEPPS pollen-vegetation model (Dawson et al., 2016; Goring et al., 2016;
906 Kujawa et al., 2016).

907 **Building Partnerships with Allied Resources**

908 In the paleosciences, a distributed network of data resources has emerged, each serving a
909 particular suite of data and research communities: evolutionary biologists, paleoclimatologists,
910 archaeologists, sample and core curators, etc. From this perspective, Neotoma is one node
911 among several in what is emerging as a federated ecosystem of complementary and allied data

912 resources. The key need is to interlink these resources through adoption of common standards
913 and data identifiers, so that scientific users can easily gather data from multiple sources.

914 Initial efforts to interlink resources have been supported by the NSF EarthCube program
915 and include the Cyberinfrastructure for Paleogeoscience (C4P) research coordination network
916 (RCN) and the EarthRates RCN, which have brought together leaders and users of
917 cyberinfrastructure resources in the paleogeosciences to identify priorities for developing
918 common standards and integrative analytical tools. Neotoma belongs to the EarthCube Council
919 of Data Facilities (<https://earthcube.org/group/council-data-facilities>). One outcome of these
920 collaborations is the Earth Life Consortium (<http://earthlifeconsortium.org/>), which is building
921 easy-to-adopt APIs that can simultaneously search for data from multiple paleobiological data
922 repositories. We are also working with VertNet (<http://vertnet.org/>) to send Neotoma data to the
923 Global Biodiversity Information Facility (GBIF, <http://www.gbif.org/>).

924 **Adding Data to Neotoma**

925 Entering, preparing, and validating data for entry into Neotoma requires effort. This effort of
926 data cleaning, validation, and preparation is not unique to Neotoma, of course, and it remains
927 the single largest bottleneck to comprehensive global-scale syntheses of paleoenvironmental
928 data. A great advantage of putting data into Neotoma is that this effort is done just once, and
929 then the data are readily available to multiple groups for multiple data synthesis projects. This
930 unified effort contrasts with the common ad hoc and inefficient system in which multiple groups
931 search for, obtain, and clean the same individual datasets.

932 Multiple solutions exist to this data ingest barrier, and Neotoma is exploring them all.
933 One ready-to-go solution is to crowdsource it, by encouraging and incentivizing individual
934 scientists to contribute data (e.g. through training workshops, recognition of Data Stewards, and
935 by supporting scientists' research objectives). Another solution is to further streamline data
936 input systems that prepare data and metadata in import-ready formats for Neotoma, at the
937 earliest stages of scientific workflows (e.g. building customized import software for data
938 generated from mass spectrometers or for microfossil analysts working at microscopes). A third
939 is to make use of text-mining programs such as GeoDeepDive (Peters et al., 2017; Peters et al.,
940 2014) for help discover and bring on-line dark data (Heidorn, 2008) and automating the
941 discovery of new papers and the particularly routine parts of data entry, such as bibliographic
942 citations. We welcome innovation by others in all of these areas.

943 **Scalability and Sustainability**

944 Scalability refers to the ability of Neotoma to grow, both with respect to the size of its
945 data holdings and community of scientific data contributors and users, while sustainability refers
946 to the ability of Neotoma to persist over time. The two are closely linked, because in both cases
947 the central solution rests in Neotoma as resource that both supports community research
948 priorities, and is supported by its constituent communities.

949 With respect to scalability, the key barriers are primarily social, rather than technological.
950 In particular, potential limits to scalability include Neotoma's emphasis on 1) expert curation to
951 ensure high data quality and utility and 2) Neotoma's model of community governance and
952 voluntary participation. The former requires both community crowdsourcing and streamlined
953 support tools for data entry, validation, and correction, while the latter means that Neotoma
954 adoption is a voluntary social process, which can be fast or slow. Neotoma is working on the
955 scalability challenge through the mechanisms described in the *Governance Model* and *New*
956 *Users, Contributors, and Communities* sections above.

957 Sustainability covers multiple dimensions, including guarantees of the long-term
958 preservation of the data housed by Neotoma and sustainability of the development efforts linked
959 to Neotoma. Data sustainability is the easier challenge: Neotoma ensures long-term data
960 sustainability through multiple and redundant mechanisms, including partnerships with multiple
961 organizations and by ensuring that Neotoma and its constituent databases exist in multiple
962 repositories (*Technical Specifications and Software Ecosystem*). Data sustainability solutions
963 are well developed.

964 Sustainability of development is the larger challenge. Neotoma development efforts so
965 far have been primarily supported by the US National Science Foundation. This reliance on a
966 single funding source is risky, but no clear alternatives yet exist. Ongoing hosting of Neotoma
967 data and data services is cheap, requiring just a few servers, and could be maintained
968 indefinitely, in case of a lapse in funding. The real risk involves sustaining the community of
969 developers who build and update these services. Neotoma's data are complex, and its
970 infrastructure development relies on individuals with good practical training in both data
971 sciences and paleosciences. This talent pool is very small but growing, and recruitment and
972 retention of this talent remains a persistent challenge.

973 Community engagement is critical to sustainability because if Neotoma is closely engaged
974 with its data contributors and users, and is seen by the scientific community as vital
975 cyberinfrastructure that facilitates large-scale earth system science, then its prospects improve

976 for continued support, growth, and development. If Neotoma is not serving this mission, then it
977 will (and should) ultimately lapse. Our personal view is that funding agencies need to commit
978 resources to long-term support of meso-scale community cyberinfrastructure efforts such as
979 Neotoma, as they do to other forms of physical scientific infrastructure, contingent upon
980 satisfactory demonstration that these efforts are advancing community scientific needs.

981 One charge for the Neotoma Leadership Council is to explore multiple funding sources
982 and business models. Because of Neotoma's commitment to open data, there are no plans to
983 charge users for data access, but other options exist. One is for other national science
984 agencies to support Neotoma's data ingest and development activities, perhaps through the
985 leadership of individual investigators associated with Constituent Databases or new data
986 synthesis efforts. A second is to establish a voluntary dues model, perhaps through partnership
987 with professional societies, which traditionally have supported other forms of scientific
988 knowledge dissemination such as peer-reviewed journals. A third is to partner with journals by
989 providing them with a high-quality data archival service that meets community data standards.
990 A fourth is to partner with investigators on data mobilization campaigns that as part of the
991 process include resources for preparing and vetting datasets for upload to Neotoma. A fifth is to
992 work with home universities to establish long-term base support for scientific databases, in line
993 with universities' mission of discovering and disseminating knowledge, in a role similar to that
994 served by university presses. All these options are viable, and all are being explored.

995 While recognizing that uncertainties exist, we are fundamentally optimistic about the long-
996 term persistence, growth, and evolution of Neotoma. Neotoma, through its constituent
997 databases, has supported macro-scale research for decades, and it has coalesced and grown
998 organically from the lab-scale data synthesis efforts of individual investigators, to the
999 development of relational database systems, and now to the development of online client-server
1000 architectures and the rise of distributed and networked networks of developers and scientists.
1001 Neotoma originated in direct response to the scientific objectives of paleoecologists and allied
1002 disciplines, and the general challenge of pursuing broad-scale science with local-scale data. As
1003 data volumes grow, both inside Neotoma (Fig. 6) and outside, community-curated data
1004 resources such as Neotoma are, increasingly, foundational infrastructure for big data science.

1005 CONCLUSIONS

1006 The Neotoma Paleoecology Database seeks to advance large-scale paleoecological,
1007 biogeographic, and global change research by providing an open, high-quality, and community-

1008 curated resource for paleoecological and associated paleoenvironmental data. Sedimentary
1009 paleoecological proxy data are expensive to collect, in time and money; Neotoma provides a
1010 low-cost solution to data sharing and access via a common platform for many different kinds of
1011 paleoecological and associated data. High data quality is achieved through open and
1012 distributed scientific governance, based on a distributed network of expert Data Stewards and
1013 associated Constituent Databases. Neotoma is in a growth stage, with open doors for
1014 membership, new Data Stewards being trained, development of new functional capabilities,
1015 extension of the data model to additional data types (e.g. organic biomarkers, stable isotopes),
1016 and more data uploaded. At the same time, much more work remains to be done, given the
1017 large volumes of paleoecological data worldwide that remain dark, trapped in unstructured
1018 publication supplements, spreadsheets on personal computers, or other inaccessible venues,
1019 and at high risk of permanent loss. We respectfully encourage other paleoecologists,
1020 paleontologists, paleoclimatologists, archaeologists, and allied disciplines to use Neotoma data
1021 and software resources as part of their research workflows, to contribute their paleoecological
1022 and associated data to Neotoma (or other community-curated resources), and to serve as
1023 members, Data Stewards, Taxonomic Experts, and on the Leadership Council. Such service
1024 can advance both personal research goals for one's own region, time, taxonomic group, and
1025 questions of interest and broader community goals of open data and enabling large-scale
1026 science. Gathering, structuring, and sharing our hard-won data into larger open resources is
1027 our community's big data challenge; community-curated resources such as Neotoma are an
1028 essential part of our community's solution.

1029 **ACKNOWLEDGMENTS**

1030 Neotoma is a community effort, and it relies on the voluntary data contributions made by
1031 individual scientists and research groups and on work by Data Stewards (<http://bit.ly/2tzjEsZ>) to
1032 clean and check these data during and after their submission to Neotoma. We thank them all.
1033 Neotoma has been supported by the NSF Geoinformatics (0948652, 0947459, 1550707,
1034 1550717, 1550805, 1550728, 1550716, 1550700, 1550890, 1550721, 1550755) and EarthCube
1035 (1541002, 1540994, 1541015, 1540979, 1540977) programs and the Wisconsin Alumni
1036 Research Foundation (WARF). Additional support came from the Wisconsin Alumni Research
1037 Foundation. Any use of trade, firm, or product names is for descriptive purposes and does not
1038 imply endorsement by the US government. Much of this paper was written while JWW was a
1039 visiting fellow at Durham University, hosted by the Institute for Advanced Study and Dr. Brian

1040 Huntley. This manuscript was improved by comments from Tom Webb, Simon Brewer, and
1041 several anonymous reviewers. Scott Farley assisted in figure design and drafting. Neotoma is
1042 indebted to the vision and work of early builders of community paleodatabases, including Pat
1043 Anderson, Anthony Barnosky, Patrick Bartlein, Richard Bradshaw, Simon Brewer, Paul
1044 Buckland, Rachid Cheddadi, Jacques-Louis de Beaulieu, Joel Guiot, Sheila Hicks, Geoff Hope,
1045 Brian Huntley, Anne-Marie Lézine, Anatoly Lozhkin, Ernie Lundelius, Vera Markgraf, Pierre
1046 Richard, Jon Sadler, and Tom Webb, as well as to all the database coordinators and
1047 contributors too numerous to mention. Their vision carries on today.

1048 TABLES

1049 **Table 1:** Constituent Databases in Neotoma and the number of datasets in each.

1050 FIGURES

1051
1052 **Figure 1:** Papers citing Neotoma and its constituent databases.
1053
1054 **Figure 2:** Neotoma serves many communities and acts as a boundary organization (Guston,
1055 2001) among these communities. Neotoma serves paleoecologists by providing a high-quality
1056 repository for their paleoecological data, with value added via digital object identifiers to facilitate
1057 data citation, data curation, and a flexible data model. Neotoma serves data users by providing
1058 a well-structured, open-access, and easy-to-use source of paleoecological data, specializing in
1059 timescales that bridge the boundary between global change ecology and geology (Betancourt,
1060 2012; Dietl and Flessa, 2011; Jackson, in press; Jackson and Blois, 2015; Jackson and Hobbs,
1061 2009; Kidwell, 2015). In return, these communities generate new questions and analytical
1062 approaches for paleoecological data. Neotoma serves educators, students, and the general
1063 public seeking to learn about the past distributions of charismatic species such as the
1064 Pleistocene megafauna and the effects of climate change on species distribution and diversity.
1065 Neotoma also serves as a boundary organization between geoscientists and computer
1066 scientists, passing data, new research questions, best practices and protocols, and geoscientific
1067 use cases & priorities.

1068

1069 **Figure 3:** Diagram of Neotoma's governance structure. Neotoma is governed by a Leadership
1070 Council, which is populated by elected members serving four year terms. The Executive
1071 working group coordinates day-to-day operations and reports to the Leadership Council. Other
1072 working groups coordinate education and outreach activities, build informatics and development
1073 activities, cultivate international partnerships, and handle membership requests and leadership
1074 elections. Constituent Databases and the Data Stewards within these databases are charged
1075 with uploading data to Neotoma, setting data standards and vocabularies, adopting and
1076 harmonizing taxonomies, and deciding default age models. These constituent databases are
1077 organized by taxonomic group or paleoecological proxy type and often are further subdivided by
1078 region or time period. The Neotoma governance system is extensible, such that new members
1079 can readily join and new Constituent Databases can form.

1080

1081 **Figure 4:** Diagram of the Neotoma software ecosystem. Data preparation and cleaning for
1082 upload to Neotoma are handled by the Tilia software (<https://www.tiliait.com/>), which has
1083 password-protected access for Data Stewards to upload datasets, update age models, and
1084 correct errors. Data are stored in the Neotoma relational database, which is deployed in SQL
1085 Server and currently hosted at Penn State's Center for Environmental Informatics. Neotoma
1086 data can be discovered, explored, viewed, and obtained through multiple platforms. Neotoma
1087 Explorer and its graphical map-based interface is designed for first-pass data explorations, new
1088 users, and educational and student groups. The APIs and neotoma R package are intended for
1089 programmatic access and for users who wish to do large-volume searches of Neotoma data
1090 holdings. Tilia can also download datasets from Neotoma, which is useful for data
1091 visualizations and for Data Stewards needing to update datasets or looking for examples of
1092 prepared Tilia files.

1093

1094 **Figure 5:** The Neotoma data model handles different kinds of sampling designs by
1095 paleoecologists through a flexible hierarchical system consisting of sites, collection units,
1096 analysis units, samples, and datasets. **Sites** are the field locations from which paleoecological
1097 data are obtained and can contain multiple collection units. **Collection units** are the specific
1098 point-level locations within sites from which data are obtained and can contain multiple analysis
1099 units. **Analysis units** are the specific depth horizons from which data are obtained and can
1100 contain multiple samples. A **sample** is a single piece of material extracted from an analysis
1101 unit, for which a single kind of measurement is made (e.g. analyzed for fossil pollen, stable

1102 isotopic analyses, etc.). A **dataset** comprises all samples of a single dataset type in a single
1103 collection unit, e.g. all pollen samples from a single core.

1104

1105 **Figure 6:** History of data uploads to Neotoma, expressed as number of datasets (a) and
1106 observations (b). Neotoma launched in 2009 with a number of datasets already in it, mostly
1107 pollen and vertebrates, representing prior database building efforts from the Global Pollen
1108 Database and FAUNMAP efforts. Rate of data uploads accelerated after 2013, when the new
1109 Neotoma data model was established and Tilia's data upload and validation routines were
1110 written. The number of datasets is relatively even among several major dataset types
1111 (vertebrates, pollen, geochronological data) with recent rapid growth of ostracode and diatom
1112 datasets. The number of pollen observations (b) is large relative to the number of datasets (a)
1113 because pollen datasets often have many samples (e.g. many samples per core) and many
1114 variables per sample (i.e. dozens of taxa per sample). As other taxa- and sample-rich datasets
1115 are added to Neotoma (e.g. diatoms, ostracodes), their relative proportion will quickly increase.

1116

1117 **Figure 7:** Tilia's interface for stewards to add new taxonomic names to Neotoma's Taxa table.
1118 Names are placed within a taxonomic tree and each taxon name is assigned a unique identifier.
1119 Stewards can also upload a citation for the source of that taxonomic name.

1120

1121

1122 **REFERENCES**

- 1123 Ammann, B., van Leeuwen, J.F.N., van der Knaap, W.O., Lischke, H., Heiri, O., Tinner, W.,
1124 2013. Vegetation responses to rapid warming and to minor climatic fluctuations during the
1125 Late-Glacial Interstadial (GI-1) at Gerzensee (Switzerland). *Palaeogeography, Palaeoclimatology, Palaeoecology* 391, Part B, 40-59.
- 1126 Arroyo-Cabral, J., Polaco, O.J., Johnson, E., 2007. An overview of the Quaternary mammals
1127 of Mexico. *Courier Forschungsinstitut Senckenberg* 259, 191-203.
- 1128 Arroyo-Cabral, J., Polaco, O.J., Johnson, E., 2009. Providing a national perspective on
1129 Quaternary mammals through a Mexican database. *The SAA Archaeological Records* 9, 21-
1130 23.
- 1131 Barnosky, A.D., Hadly, E.A., Gonzalez, P., Head, J., Polly, P.D., Lawing, A.M., Eronen, J.T.,
1132 Ackerly, D.D., Alex, K., Biber, E., Blois, J., Brashares, J., Ceballos, G., Davis, E., Dietl, G.P.,
1133 Dirzo, R., Doremus, H., Fortelius, M., Greene, H.W., Hellmann, J., Hickler, T., Jackson, S.T.,
1134 Kemp, M., Koch, P.L., Kremen, C., Lindsey, E.L., Looy, C., Marshall, C.R., Mendenhall, C.,
1135 Mulch, A., Mychajliw, A.M., Nowak, C., Ramakrishnan, U., Schnitzler, J., Das Shrestha, K.,
1136 Solari, K., Stegner, L., Stegner, M.A., Stenseth, N.C., Wake, M.H., Zhang, Z., 2017. Merging
1137 paleobiology with conservation biology to guide the future of terrestrial ecosystems. *Science*
1138 355.
- 1139 Barnosky, A.D., Lindsey, E.L., Villavicencio, N.A., Bostelmann, E., Hadly, E.A., Wanket, J.,
1140 Marshall, C.R., 2016. Variable impact of late-Quaternary megafaunal extinction in causing
1141 ecological state shifts in North and South America. *Proceedings of the National Academy of
1142 Sciences* 113, 856-861.
- 1143 Bartlein, P.J., Harrison, S.P., Brewer, S., Connor, S., Davis, B.A.S., Gajewski, K., Guiot, J.,
1144 Harrison-Prentice, T.I., Henderson, A., Peyron, O., Prentice, I.C., Scholze, M., Seppä, H.,
1145 Shuman, B., Sugita, S., Thompson, R.S., Viau, A.E., Williams, J., Wu, H., 2011. Pollen-
1146 based continental climate reconstructions at 6 and 21 ka: a global synthesis. *Clim. Dyn.* 37,
1147 775-802.
- 1148 Bennett, C.R., Provan, J., 2008. What do we mean by 'refugia'? *Quaternary Science Reviews*
1149 27, 2449-2455.
- 1150 Bernabo, J.C., Webb, T., III, 1977. Changing patterns in the Holocene pollen record of
1151 northeastern North America: A mapped summary. *Quaternary Research* 8, 64-96.
- 1152 Betancourt, J., 2012. Reflections on the relevance of history in a nonstationary world, in: Wiens,
1153 J., Hayward, G.D., Safford, H.D., Giffen, C.M. (Eds.), *Historical Environmental Variation in*

- 1155 Conservation and Natural Resource Management, First edition. John Wiley & Sons, Ltd., pp.
1156 307-317.
- 1157 Betancourt, J.L., Van Devender, T.R., Martin, P.S., 1990. *Packrat Middens: The Last 40,000*
1158 Years of Biotic Change
- Tucson, The University of Arizona Press, p. 467.
- 1159 Birks, H.H., 2015. South to north: Contrasting late-glacial and early-Holocene climate changes
1160 and vegetation responses between south and north Norway. *Holocene* 25, 37-52.
- 1161 Birks, H.J.B., 1995. Quantitative paleoenvironmental reconstruction, in: Maddy, D., Brew, J.S.
1162 (Eds.), *Statistical modelling of Quaternary Science Data*. Technical Guide 5. Quaternary
1163 Research Association, Cambridge, pp. 116-254.
- 1164 Birks, H.J.B., Line, J.M., 1992. The use of rarefaction analysis for estimating palynological
1165 richness from Quaternary pollen-analytical data. *The Holocene* 2, 1-10.
- 1166 Blaauw, M., 2010. Methods and code for 'classical' age-modelling of radiocarbon sequences.
1167 *Quaternary Geochronology* 5, 512-518.
- 1168 Blaauw, M., Christen, J.A., 2011. Flexible paleoclimate age-depth models using an
1169 autoregressive gamma process. *Bayesian Analysis* 6, 1-18.
- 1170 Blarquez, O., Carcailliet, C., Frejaville, T., Bergeron, Y., 2014. Disentangling the trajectories of
1171 alpha, beta and gamma plant diversity of North American boreal ecoregions since 15,500
1172 years. *Frontiers in Ecology and Evolution* 2.
- 1173 Blois, J.L., Williams, J.W., Grimm, E.C., Jackson, S.T., Graham, R.W., 2011. A methodological
1174 framework for improved paleovegetation mapping from late-Quaternary pollen records.
1175 *Quaternary Science Reviews* 30, 1926-1939.
- 1176 Blois, J.L., Zarnetske, P.L., Fitzpatrick, M.C., Finnegan, S., 2013. Climate change and the past,
1177 present, and future of biotic interactions. *Science* 341, 499-504.
- 1178 Bocinsky, R.K., Kohler, T.A., 2014. A 2,000-year reconstruction of the rain-fed maize agricultural
1179 niche in the US Southwest. *Nature Communications* 5, 5618.
- 1180 Boyle, J., 1999. Variability of tephra in lake and catchment sediments, Svínnavatn, Iceland.
1181 *Global and Planetary Change* 21, 129-149.
- 1182 Brewer, S., Jackson, S.T., Williams, J.W., 2012. Paleoeconinformatics: Applying geohistorical
1183 data to ecological questions. *Trends Ecol. Evol.* 27, 104-112.
- 1184 Buckland, P.I., 2007. The Development and Implementation of Software for
1185 Palaeoenvironmental and Palaeoclimatological Research: The Bugs Coleopteran Ecology
1186 Package (BugsCEP), Department of Archaeology & Sámi Studies. University of Umeå,
1187 Umeå, Sweden, p. 236 pp + CD.

- 1188 Bush, R.T., McInerney, F.A., 2013. Leaf wax n-alkane distributions in and across modern plants:
1189 Implications for paleoecology and chemotaxonomy. *Geochimica et Cosmochimica Acta* 117,
1190 161-179.
- 1191 Carrasco, M.A., Barnosky, A.D., Kraatz, B.P., Davis, E.B., 2007. The Miocene mammal
1192 mapping project (MIOMAP): An online database of Arikareean through Hemphillian fossil
1193 mammals. *Bulletin of Carnegie Museum of Natural History* 39, 183-188.
- 1194 Cinget, B., de Lafontaine, G., Gérardi, S., Bousquet, J., 2015. Integrating phylogeography and
1195 paleoecology to investigate the origin and dynamics of hybrid zones: insights from two
1196 widespread North American firs. *Molecular Ecology* 24, 2856-2870.
- 1197 Clarke, S.J., Lynch, A.J.J., 2016. Palaeoecology to inform wetland conservation and
1198 management: some experiences and prospects. *Marine and Freshwater Research*
1199 <http://dx.doi.org/10.1071/MF15031>.
- 1200 Clement, B.M., 2004. Dependence of the duration of geomagnetic polarity reversals on site
1201 latitude. *Nature* 428, 637-640.
- 1202 CLIMAP Project Members, 1976. The surface of Ice-Age Earth: Quantitative geological
1203 evidence is used to reconstruct boundary conditions for the climate 18,000 years ago.
1204 *Science* 191, 1131-1137.
- 1205 Davis, M.B., 1976. Pleistocene biogeography of temperate deciduous forests. *Geoscience and*
1206 *Man* XIII, 13-26.
- 1207 Dawson, A., Paciorek, C.J., McLachlan, J.S., Goring, S., Williams, J.W., Jackson, S.T., 2016.
1208 Quantifying pollen-vegetation relationships to reconstruct forests using 19th-century forest
1209 composition and pollen data. *Quaternary Science Reviews* 137, 156-175.
- 1210 Dawson, T.P., Jackson, S.T., House, J.I., Prentice, I.C., Mace, G.M., 2011. Beyond predictions:
1211 Biodiversity conservation in a changing climate. *Science* 332, 53-58.
- 1212 De La Torre, A.R., Roberts, D.R., Aitken, S.N., 2014. Genome-wide admixture and ecological
1213 niche modelling reveal the maintenance of species boundaries despite long history of
1214 interspecific gene flow. *Molecular Ecology* 23, 2046-2059.
- 1215 deMenocal, P.B., 2001. Cultural responses to climate change during the late Holocene. *Science*
1216 292, 667-673.
- 1217 Dietl, G.P., Flessa, K.W., 2011. Conservation paleoecology: Putting the dead to work. *Trends*
1218 *Ecol. Evol.* 26, 30-37.
- 1219 Dietl, G.P., Kidwell, S.M., Brenner, M., Burney, D.A., Flessa, K.W., Jackson, S.T., Koch, P.L.,
1220 2015. Conservation paleobiology: Leveraging knowledge of the past to inform conservation
1221 and restoration. *Annual Review of Earth and Planetary Sciences* 43, 79-103.

- 1222 Dietze, M.C., 2017. Ecological Forecasting. Princeton University Press, Princeton, NJ.
- 1223 Dietze, M.C., Lebauer, D.S., Kooper, R., 2012. On improving the communication between
- 1224 models and data. *Plant, Cell & Environment* doi: 10.1111/pce.12043.
- 1225 Doughty, C.E., Wolf, A., Madhi, Y., 2013. The legacy of the Pleistocene megafauna extinctions
- 1226 on nutrient availability in Amazonia. *Nature Geoscience* 6, 761-764.
- 1227 ECMA International, 2013. The JSON data standard format. ECMA International, Geneva,
- 1228 Switzerland.
- 1229 Ellis, E.C., Kaplan, J.O., Fuller, D.Q., Vavrus, S., Goldewijk, K.K., Verburg, P.H., 2013. Used
- 1230 planet: A global history. *Proceedings of the National Academy of Sciences* 110, 7978-7985.
- 1231 Emery-Wetherell, M.M., McHorse, B.K., Davis, E.B., in press. Spatially explicit analysis sheds
- 1232 new light on the Pleistocene Megafaunal Extinction in North America. *Paleobiology*.
- 1233 Evans, M.N., Tolwinski-Ward, S.E., Thompson, D.M., Anchukaitis, K.J., 2013. Applications of
- 1234 proxy system modeling in high resolution paleoclimatology. *Quaternary Science Reviews*
- 1235 76, 16-28.
- 1236 FAUNMAP Working Group, 1994. FAUNMAP: A Database Documenting Late Quaternary
- 1237 Distributions of Mammal Species in the United States. Illinois State Museum, Springfield, IL.
- 1238 Ferguson, A.R., Nielson, J.L., Cragin, M.H., Bandrowski, A.E., Martone, M.E., 2014. Big data
- 1239 from small data: data-sharing in the 'long tail' of neuroscience. *Nature Neuroscience* 17,
- 1240 1442-1448.
- 1241 Finsinger, W., Giesecke, T., Brewer, S., Leydet, M., 2017. Emergence patterns of novelty in
- 1242 European vegetation assemblages over the past 15 000 years. *Ecology Letters* 20, 336–
- 1243 346.
- 1244 Flantua, S.G.A., Hooghiemstra, H., Grimm, E.C., Behling, H., Bush, M.B., González-Arango, C.,
- 1245 Gosling, W.D., Ledru, M.-P., Lozano-García, S., Maldonado, A., Prieto, A.R., Rull, V., Van
- 1246 Boxel, J.H., 2015. Updated site compilation of the Latin American Pollen Database. *Rev.*
- 1247 *Palaeobot. Palynology* 223, 104-115.
- 1248 Forester, R.M., Smith, A.J., Palmer, D.F., Curry, B.B., 2005. NANODe: North American
- 1249 Nonmarine Ostracode Database, version 1. Kent State University, Kent, OH, USA.
- 1250 Fritz, S.A., Schnitzler, J., Eronen, J.T., Hof, C., Böhning-Gaese, K., Graham, C.H., 2013.
- 1251 Diversity in time and space: wanted dead and alive. *Trends Ecol. Evol.* 28, 509-516.
- 1252 Fyfe, R.M., de Beaulieu, J.-L., Binney, H., Bradshaw, R.H.W., Brewer, S., Le Flao, A., Finsinger,
- 1253 W., Gaillard, M.-J., Giesecke, T., Gil-Romera, G., Grimm, E.C., Huntley, B., Kunes, P., Kühl,
- 1254 N., Leydet, M., Lotter, A.F., Tarasov, P.E., Tonkov, S., 2009. The European Pollen
- 1255 Database: past efforts and current activities. *Veg. Hist. Archaeobot.* 18, 417-424.

- 1256 Gaillard, M.J., Sugita, S., Mazier, F., Trondman, A.K., Broström, A., Hickler, T., Kaplan, J.O.,
1257 Kjellström, E., Kokfelt, U., Kuneš, P., Lemmen, C., Miller, P., Olofsson, J., Poska, A.,
1258 Rundgren, M., Smith, B., Strandberg, G., Fyfe, R., Nielsen, A.B., Alenius, T., Balakauskas,
1259 L., Barnekow, L., Birks, H.J.B., Bjune, A., Björkman, L., Giesecke, T., Hjelle, K., Kalnina, L.,
1260 Kangur, M., van der Knaap, W.O., Koff, T., Lagerås, P., Latałowa, M., Leydet, M.,
1261 Lechterbeck, J., Lindbladh, M., Odgaard, B., Peglar, S., Segerström, U., von Stedingk, H.,
1262 Seppä, H., 2010. Holocene land-cover reconstructions for studies on land cover-climate
1263 feedbacks. *Climates of the Past* 6, 483-499.
- 1264 Gavin, D.G., Fitzpatrick, M.C., Gugger, P.F., Heath, K.D., Rodríguez-Sánchez, F., Dobrowski,
1265 S.Z., Hampe, A., Hu, F.S., Ashcroft, M.B., Bartlein, P.J., Blois, J.L., Carstens, B.C., Davis,
1266 E.B., de Lafontaine, G., Edwards, M.E., Fernandez, M., Henne, P.D., Herring, E.M., Holden,
1267 Z.A., Kong, W.-s., Liu, J., Magri, D., Matzke, N.J., McGlone, M.S., Saltré, F., Stigall, A.L.,
1268 Tsai, Y.-H.E., Williams, J.W., 2014. Climate refugia: joint inference from fossil records,
1269 species distribution models and phyogeography. *New Phytologist* 204, 37-54.
- 1270 Giesecke, T., Brewer, S., Finsinger, W., Leydet, M., Bradshaw, R.H.W., 2017. Patterns and
1271 dynamics of European vegetation change over the last 15,000 years. *Journal of*
1272 *Biogeography* 44, 1441-1456.
- 1273 Giesecke, T., Davis, B., Brewer, S., Finsinger, W., Wolters, S., Blaauw, M., de Beaulieu, J.-L.,
1274 Binney, H., Fyfe, R., Gaillard, M.-J., Gil-Romera, G., van der Knaap, W.O., Kuneš, P., Kühl,
1275 N., van Leeuwen, J.N., Leydet, M., Lotter, A., Ortú, E., Semmler, M., Bradshaw, R.W., 2014.
1276 Towards mapping the late Quaternary vegetation change of Europe. *Veg. Hist. Archaeobot.*
1277 23, 75-86.
- 1278 Goring, S., Dawson, A., Simpson, G., Ram, K., Graham, R.W., Grimm, E.C., Williams, J.W.,
1279 2015. *neotoma*: A Programmatic Interface to the Neotoma Paleoecological Database. *Open*
1280 *Quaternary* 1, 1-17.
- 1281 Goring, S.J., Williams, J.W., 2017. Effect of historic land-use and climate change on tree-climate
1282 relationships in the upper Midwestern United States. *Ecology Letters*.
- 1283 Goring, S.J., Williams, J.W., Mladenoff, D.J., Cogbill, C.V., Record, S., Paciorek, C.J., Jackson,
1284 S.J., Dietze, M.C., McLachlan, J.S., 2016. Novel and lost forests in the upper Midwestern
1285 United States, from new estimates of settlement-era composition, stem density, and
1286 biomass. *PLoS One* 11, e0151935.
- 1287 Graham, R.W., Lundelius, E.L., Jr., Graham, M.A., Schroeder, E.K., Toomey, R.S., III,
1288 Anderson, E., Barnosky, A.D., Burns, J.A., Churcher, C.S., Grayson, D.K., Guthrie, R.D.,
1289 Harrington, C.R., Jefferson, G.T., Martin, L.D., McDonald, H.G., Morlan, R.E., Semken Jr.,

- 1290 H.A., Webb, S.D., Werdelin, L., Wilson, M.C., 1996. Spatial response of mammals to Late
1291 Quaternary environmental fluctuations. *Science* 272, 1601:1606.
- 1292 Grant, M.J., Stevens, C.J., Whitehouse, N.J., Norcott, D., Macphail, R.I., Langdon, C.,
1293 Cameron, N., Barnett, C., Langdon, P.G., Crowder, J., Mulhall, N., Attree, K., Leivers, M.,
1294 Greatorex, R., Ellis, C., 2014. A palaeoenvironmental context for Terminal Upper
1295 Palaeolithic and Mesolithic activity in the Colne Valley: Offsite records contemporary with
1296 occupation at Three Ways Wharf, Uxbridge. *Environmental Archaeology* 19, 131-152.
- 1297 Grimm, E.C., 1988. Data analysis and display, in: Huntley, B., Webb, T., III (Eds.), *Vegetation*
1298 *History*. Kluwer Academic Publishers, Dordrecht, pp. 43-76.
- 1299 Grimm, E.C., Blaauw, M., Buck, C.E., Williams, J.W., 2014. Age models, chronologies, and
1300 databases workshop: Complete report and recommendations. *PAGES Workshop Report* 22.
- 1301 Grimm, E.C., Keltner, J., Cheddadi, R., Hicks, S., Lézine, A.-M., Berrio, J.C., Williams, J.W.,
1302 2013. Pollen databases and their application, in: Elias, S.A., Mock, C.J. (Eds.),
1303 *Encyclopedia of Quaternary Science*. Elsevier, pp. 831-838.
- 1304 Guston, D.H., 2001. Boundary organizations in environmental policy and science: An
1305 introduction. *Science, Technology, & Human Values* 26, 399-408.
- 1306 Gutiérrez-García, T.A., Vázquez-Domínguez, E., Arroyo-Cabralles, J., Kuch, M., Enk, J., King,
1307 C., Poinar, H.N., 2014. Ancient DNA and the tropics: a rodent's tale. *Biology Letters* 10,
1308 20140224.
- 1309 Hampton, S.E., Strasser, C.A., Tewksbury, J.J., Gram, W.K., Budden, A.E., Batcheller, A.L.,
1310 Duke, C.S., Porter, J.H., 2013. Big data and the future of ecology. *Frontiers in Ecology and*
1311 *the Environment* 11, 156-162.
- 1312 Heffernan, J.B., Soranno, P.A., Angilletta, M.J., Buckley, L.B., Gruner, D.S., Keitt, T.H., Kellner,
1313 J.R., Kominoski, J.S., Rocha, A.V., Xiao, J., Harms, T.K., Goring, S.J., Koenig, L.E.,
1314 McDowell, W.H., Powell, H., Richardson, A.D., Stow, C.A., Vargas, R., Weathers, K.C.,
1315 2014. Macrosystems ecology: understanding ecological patterns and processes at
1316 continental scales. *Frontiers in Ecology and the Environment* 12, 5-14.
- 1317 Heidorn, P.B., 2008. Shedding light on the dark data in the long tail of science. *Library Trends*
1318 57, 280-299.
- 1319 Howe, D., Costanzo, M., Fey, P., Gojobori, T., Hannick, L., Hide, W., Hill, D.P., Kania, R.,
1320 Schaeffer, M., St Pierre, S., Twigger, S., White, O., Yon Rhee, S., 2008. Big data: The future
1321 of biocuration. *Nature* 455, 47-50.
- 1322 Huntley, B., Birks, H.J.B., 1983. *An Atlas of Past and Present Pollen Maps for Europe: 0-13000*
1323 *Years Ago*. Cambridge University Press, Cambridge.

- 1324 Jackson, S.T., in press. Late Quaternary biogeography: linking biotic responses to
1325 environmental variability across timescales, in: Lomolino, M., Heaney, L. (Eds.), *Frontiers of*
1326 *Biogeography: New Directions in the Geography of Nature*. Sinauer Associates, Inc.,
1327 Sunderland, MA.
- 1328 Jackson, S.T., Blois, J.L., 2015. Community ecology in a changing environment. *Proceedings of*
1329 *the National Academy of Sciences* 112, 4915-4921.
- 1330 Jackson, S.T., Grimm, E.C., Thompson, R.S., 2000. Database resources in Quaternary
1331 paleobotany, in: Lipscomb, B., Pipoly, J., Sanders, R. (Eds.), *Floristics in the New*
1332 *Millennium*, pp. 113-120.
- 1333 Jackson, S.T., Hobbs, R.J., 2009. Ecological restoration in the light of ecological history.
1334 *Science* 325, 567-569.
- 1335 Jackson, S.T., Overpeck, J.T., Webb, T., III, Keatitch, S.E., Anderson, K.H., 1997. Mapped
1336 plant-macrofossil and pollen records of late Quaternary vegetation change in eastern North
1337 America. *Quaternary Science Reviews* 16, 1-70.
- 1338 Jezkova, T., Riddle, B.R., Card, D.C., Schield, D.R., Eckstut, M.E., Castoe, T.A., 2015. Genetic
1339 consequences of postglacial range expansion in two codistributed rodents (genus
1340 *Dipodomys*) depend on ecology and genetic locus. *Molecular Ecology* 24, 83-97.
- 1341 Juggins, S., 2015. *rioja*: Analysis of Quaternary Science Data, R package version (0.9-5).
- 1342 Kaplan, J.O., Krumhardt, K.M., Ellis, E.C., Ruddiman, W.F., Lemmen, C., Goldewijk, K.K., 2011.
1343 Holocene carbon emissions as a result of anthropogenic land cover change. *The Holocene*
1344 21, 775-791.
- 1345 Kaplan, J.O., Krumhardt, K.M., Zimmermann, N., 2009. The prehistoric and preindustrial
1346 deforestation of Europe. *Quaternary Science Reviews* 28, 3016-3034.
- 1347 Kidwell, S.M., 2015. Biology in the Anthropocene: Challenges and insights from young fossil
1348 records. *Proceedings of the National Academy of Sciences* 12, 4922-4929.
- 1349 Kujawa, E., Goring, S., Dawson, A., Calcote, R., Grimm, E.C., Hotchkiss, S.C., Jackson, S.T.,
1350 Lynch, E.A., McLachlan, J., St. Jacques, J.-M., Umphanowar, C., Jr., Williams, J.W., 2016.
1351 The effects of anthropogenic land cover change on pollen-vegetation relationships in the
1352 American Midwest. *Anthropocene* <http://dx.doi.org/10.1016/j.ancene.2016.09.005>.
- 1353 Latorre, C., Moreno, P.I., Grimm, E.C., 2014. 1st Workshop on paleoecological databases in
1354 South America. *PAGES Magazine* 22, 52.
- 1355 Lehnert, K., Hsu, L., 2015. The new paradigm of data publication. *Elements* 11, 368-369.

- 1356 Li, C., 2004. Dynamics of late Quaternary mammal population inferred from geostatistical study
1357 of the Faunmap database and Its implications for conservation. Universität Trier, Trier, p.
1358 172.
- 1359 Lisiecki, L.E., Raymo, M.E., 2005. A Pliocene-Pleistocene stack of 57 globally distributed
1360 benthic $d^{18}O$ records. *Paleoceanography* 20, DOI: 10.1029/2004PA001071.
- 1361 Loeffler, S., Ai, S., McEwan, R., Myrbo, A., 2015. Flyover Country: A plane ride could be to
1362 geoscience outreach what a planetarium is to astronomy outreach - the perfect venue for
1363 sharing big, awe inspiring ideas, with a view to match, American Geophysical Union, San
1364 Francisco, CA, pp. ED54A-07.
- 1365 Lorenzen, E.D., Nogues-Bravo, D., Orlando, L., Weinstock, J., Binladen, J., Marske, K.A., Ugan,
1366 A., Borregaard, M.K., Gilbert, M.T.P., Nielsen, R., Ho, S.Y.W., Goebel, T., Graf, K.E., Byers,
1367 D., Stenderup, J.T., Rasmussen, M., Campos, P.F., Leonard, J.A., Koepfli, K.-P., Froese,
1368 D., Zazula, G., Stafford, T.W., Aaris-Sorensen, K., Batra, P., Haywood, A.M., Singarayer,
1369 J.S., Valdes, P.J., Boeskorov, G., Burns, J.A., Davydov, S.P., Haile, J., Jenkins, D.L.,
1370 Kosintsev, P., Kuznetsova, T., Lai, X., Martin, L.D., McDonald, H.G., Mol, D., Meldgaard, M.,
1371 Munch, K., Stephan, E., Sablin, M., Sommer, R.S., Sipko, T., Scott, E., Suchard, M.A.,
1372 Tikhonov, A., Willerslev, R., Wayne, R.K., Cooper, A., Hofreiter, M., Sher, A., Shapiro, B.,
1373 Rahbek, C., Willerslev, E., 2011. Species-specific responses of Late Quaternary megafauna
1374 to climate and humans. *Nature* 479, 359-364.
- 1375 Lynch, C., 2008. How do your data grow? *Nature* 455, 28-29.
- 1376 Lyons, S.K., 2001. A quantitative assessment of the community structure and dynamics of
1377 Pleistocene mammals. University of Chicago, Chicago, p. 230.
- 1378 Marcott, S.A., Shakun, J.D., Clark, P.U., Mix, A.C., 2013. A reconstruction of regional and global
1379 temperature for the past 11,300 years. *Science* 339, 1198-1201.
- 1380 MARGO Project Members, 2009. Constraints on the magnitude and patterns of ocean cooling at
1381 the Last Glacial Maximum. *Nature Geosci* 2, 127-132.
- 1382 Martin, A.C., Harvey, W.J., 2017. The Global Pollen Project: a new tool for pollen identification
1383 and the dissemination of physical reference collections. *Methods in Ecology and Evolution*,
1384 n/a-n/a.
- 1385 Moritz, C., Agudo, R., 2013. The future of species under climate change: Resilience or decline?
1386 *Science* 341, 504-508.
- 1387 Muñoz, S.E., Gajewski, K., Peros, M.C., 2010. Synchronous environmental and cultural change
1388 in the prehistory of the northeastern United States. *Proceedings of the National Academy of
1389 Sciences* 107, 22008-22013.

- 1390 Nature, 2017. Not-so-open data. *Nature* 546, 327.
- 1391 Ogg, J.G., Smith, A.G., 2005. The geomagnetic polarity time scale., in: Gradstein, F.M., Ogg,
1392 J.G., Smith, A.G. (Eds.), *A Geologic Time Scale 2004*. Cambridge University Press,
1393 Cambridge, United Kingdom, pp. 63-86.
- 1394 Ohri, A., 2014. R with cloud APIs, R for Cloud Computing. Springer, New York, pp. 217-235.
- 1395 Olszewski, T.D., Kidwell, S.M., 2007. The preservational fidelity of evenness in molluscan death
1396 assemblages. *Paleobiology* 33, 1-23.
- 1397 Ordóñez, A., Williams, J.W., 2013. Climatic and biotic velocities for woody taxa distributions
1398 over the last 16 000 years in eastern North America. *Ecology Letters* 16, 773-781.
- 1399 Panagiotakopulu, E., Buchan, A.L., 2015. Present and Norse Greenlandic hayfields – Insect
1400 assemblages and human impact in southern Greenland. *The Holocene* 25, 921-931.
- 1401 Parnell, A.C., Haslett, J., Allen, J.R.M., Buck, C.E., Huntley, B., 2008. A flexible approach to
1402 assessing synchronicity of past events using Bayesian reconstructions of sedimentation
1403 history. *Quaternary Science Reviews* 27, 1872-1885.
- 1404 Parnell, A.C., Haslett, J., Sweeney, J., Doan, T.K., Allen, J.R.M., Huntley, B., 2016. Joint
1405 palaeoclimate reconstruction from pollen data via forward models and climate histories.
1406 *Quaternary Science Reviews* 151, 111-126.
- 1407 Peters, S.E., Husson, J.M., Wilcots, J., 2017. The rise and fall of stromatolites in shallow marine
1408 environments. *Geology* G38931.1.
- 1409 Peters, S.E., Zhang, C., Livny, M., Ré, C., 2014. A machine reading system for assembling
1410 synthetic paleontological databases. *PLoS ONE* 9, e113523.
- 1411 Radeloff, V.C., Williams, J.W., Bateman, B.L., Burke, K.D., Carter, S.K., Childress, E.S.,
1412 Cromwell, K., Gratton, C., Hasley, A.O., Kraemer, B.M., Latzka, A.W., Marin-Spiotta, E.,
1413 Meine, C.D., Munoz, S.E., Neeson, T.M., Pidgeon, A.M., Rissman, A.R., Rivera, R.J.,
1414 Szymanski, L.M., Usinowicz, J., 2015. The rise of novelty in ecosystems. *Ecol. Appl.* 25,
1415 2051-2068.
- 1416 Richard, P.J.H., 1995. Le couvert végétal du Québec-Labrador il y a 6000 ans BP: essai.
1417 *Geogr. Phys. Quat.* 49, 117-140.
- 1418 Roth, R.E., Donohue, R.G., Sack, C.M., Wallace, T.R., Buckingham, T.M.A., 2014. A process
1419 for keeping pace with evolving web mapping technologies. *Cartographic Perspectives* 25-
1420 52.
- 1421 Sachse, D., Billault, I., Bowen, G.J., Chikaraishi, Y., Dawson, T.E., Feakins, S.J., Freeman,
1422 K.H., Magill, C.R., McInerney, F.A., van der Meer, M.T.J., Polissar, P., Robins, R.J., Sachs,
1423 J.P., Schmidt, H.-L., Sessions, A.L., White, J.W.C., West, J.B., Kahmen, A., 2012. Molecular

- 1424 Paleohydrology: Interpreting the Hydrogen-Isotopic Composition of Lipid Biomarkers from
1425 Photosynthesizing Organisms. Annual Review of Earth and Planetary Sciences 40, 221-
1426 249.
- 1427 Sadler, J.P., Buckland, P.C., Rains, M., 1992. BUGS: an entomological database. Antenna 16,
1428 158-166.
- 1429 Sandom, C.J., Ejrnæs, R., Hansen, M.D.D., Svenning, J.-C., 2014. High herbivore density
1430 associated with vegetation diversity in interglacial ecosystems. Proceedings of the National
1431 Academy of Sciences 111, 4162-4167.
- 1432 Schmittner, A., Urban, N.M., Shakun, J.D., Mahowald, N.M., Clark, P.U., Bartlein, P.J., Mix,
1433 A.C., Rosell-Melé, A., 2011. Climate sensitivity estimated from temperature reconstructions
1434 of the Last Glacial Maximum. Science 334, 1385-1388.
- 1435 Seddon, A.W., Macias-Fauria, M., Willis, K.J., 2015. Climate and abrupt vegetation change in
1436 Northern Europe since the last deglaciation. The Holocene 25, 25-36.
- 1437 Shakun, J.D., Clark, P.U., He, F., Marcott, S.A., Mix, A.C., Liu, Z., Otto-Bliesner, B., Schmittner,
1438 A., Bard, E., 2012. Global warming preceded by increasing carbon dioxide concentrations
1439 during the last deglaciation. Nature 484, 49-54.
- 1440 Shuman, B.N., Newby, P., Donnelly, J.P., 2009. Abrupt climate change as an important agent of
1441 ecological change in the Northeast U.S. throughout the past 15,000 years. Quaternary
1442 Science Reviews 28, 1693-1709.
- 1443 Simpson, G.L., 2007. Analogue methods in palaeoecology: Using the analogue package
1444 Journal of Statistical Software 22, 1-29.
- 1445 Simpson, G.L., Oksanen, J., 2015. analogue: Analogue matching and Modern Analogue
1446 Technique transfer function models. (R package version 0.16-3) .
- 1447 Spaulding, W.G., Betancourt, J.L., Croft, L.K., Cole, K.L., 1990. Packrat middens: Their
1448 composition and methods of analysis, in: Betancourt, J.L., Van Devender, T.R., Martin, P.S.
1449 (Eds.), Packrat Middens: The Last 40,000 Years of Biotic Change. University of Arizona
1450 Press, Tucson, AZ, pp. 59-84.
- 1451 Sullivan, T.J., Charles, D.F., 1994. The feasibility and utility of a paleolimnology/paleoclimate
1452 data cooperative for North America. J. Paleolimn. 10, 265-273.
- 1453 Svenning, J.-C., Sandel, B., 2013. Disequilibrium vegetation dynamics under future climate
1454 change. American Journal of Botany 100, 1266-1286.
- 1455 Trouet, V., Diaz, H.F., Wahl, E.R., Viau, A.E., Graham, R., Graham, N., Cook, E.R., 2013. A
1456 1500-year reconstruction of annual mean temperature for temperate North America on
1457 decadal-to-multidecadal time scales. Environmental Research Letters 8, 024008.

- 1458 Uhen, M.D., Barnosky, A.D., Bills, B., Blois, J., Carrano, M.T., Carrasco, M.A., Erickson, G.M.,
1459 Eronen, J.T., Fortelius, M., Graham, R.W., Grimm, E.C., O'Leary, M.A., Mast, A., Piel, W.H.,
1460 Polly, P.D., Säilä, L.K., 2013. From card catalogs to computers: Databases in vertebrate
1461 paleontology. *Journal of Vertebrate Paleontology* 33, 13-28.
- 1462 Viau, A.E., Ladd, M., Gajewski, K., 2012. The climate of North America during the past 2000
1463 years reconstructed from pollen data. *Global and Planetary Change* 84–85, 75-83.
- 1464 Vickers, K., Buckland, P.I., 2015. Predicting island beetle faunas by their climate ranges: the
1465 tabula rasa/refugia theory in the North Atlantic. *Journal of Biogeography* 42, 2031-2048.
- 1466 Vincens, A., Lézine, A.-M., Buchet, G., Lewden, D., Le Thomas, A., 2007. African pollen
1467 database inventory of tree and shrub pollen types. *Rev. Palaeobot. Palynology* 145, 135-
1468 141.
- 1469 Webb, T., III, 1997. Spatial response of plant taxa to climate change: A palaeoecological
1470 perspective, in: Huntley, B., Cramer, W., Morgan, A.V., Prentice, H.C., Allen, J.R.M. (Eds.),
1471 *Past and Future Rapid Environmental Changes: The Spatial and Evolutionary Responses*
1472 of Terrestrial Biota. Springer-Verlag, Berlin, pp. 55-72.
- 1473 Weng, C., Hooghiemstra, H., Duivenvoorden, J.F., 2006. Challenges in estimating past plant
1474 diversity from fossil pollen data: statistical assessment, problems, and possible solutions.
1475 *Divers. Distrib.* 12, 310-318.
- 1476 Whitehouse, N.J., P.G., L., Bustin, R., Galsworthy, S., 2008. Fossil insects and ecosystem
1477 dynamics in wetlands: implications for biodiversity and conservation. *Biodivers. Conserv.* 17,
1478 2055–2078.
- 1479 Williams, J.W., Shuman, B.N., Webb, T., III, Bartlein, P.J., Leduc, P.L., 2004. Late Quaternary
1480 vegetation dynamics in North America: Scaling from taxa to biomes. *Ecological*
1481 *Monographs* 74, 309-334.
- 1482 Wright, H.E., Jr. , Kutzbach, J.E., Webb, T., III, Ruddiman, W.F., Street-Perrott, F.A., Bartlein,
1483 P.J., 1993. *Global Climates since the Last Glacial Maximum*. University of Minnesota Press,
1484 Minneapolis.
- 1485 Zazula, G.D., MacPhee, R.D.E., Metcalfe, J.Z., Reyes, A.V., Brock, F., Druckenmiller, P.S.,
1486 Groves, P., Harrington, C.R., Hodgins, G.W.L., Kunz, M.L., Longstaffe, F.J., Mann, D.H.,
1487 McDonald, H.G., Nalawade-Chavan, S., Southon, J.R., 2014. American mastodon
1488 extirpation in the Arctic and Subarctic predates human colonization and terminal Pleistocene
1489 climate change. *Proceedings of the National Academy of Sciences* 111, 18460-18465.
- 1490 Zdanowicz, C.M., Zielinski, G.A., Germani, M.S., 1999. Mount Mazama eruption: calendrical
1491 age verified and atmospheric impact assessed. *Geology* 27, 621-624.

1492 Zhang, Z., Zhao, M., Eglinton, G., Lu, H., Huang, C.-Y., 2006. Leaf wax lipids as
1493 paleovegetational and paleoenvironmental proxies for the Chinese Loess Plateau over the
1494 last 170kyr. Quaternary Science Reviews 25, 575-594.
1495

1496 **SUPPLEMENTARY MATERIAL**

1497 **Training and Help Resources**

1498 Neotoma PIs, Stewards, and developers have built and maintain a variety of resources
1499 designed to help orient new users, support educational and outreach activities, and to report
1500 issues. These are currently somewhat dispersed; in the future, we aim to better consolidate
1501 these resources. See Table S1 below for a full listing of Neotoma online resources.

1502 **Training Workshops.** We run several kinds of in-person training workshops and
1503 webinars. Some workshops are designed to support Neotoma data users and describe the
1504 tools available for exploring, viewing, and obtaining Neotoma data via Neotoma Explorer and
1505 the R *neotoma* package. Recent user-oriented workshops were held at the International
1506 Biogeography Society (2017), the American Quaternary Association (2016), the Society for
1507 Vertebrate Paleontology (2016), the European Pollen Database meeting (Giesecke et al. 2016),
1508 at PalEON Boot Camp (2014, 2016), the International Symposium on Ostracoda (2017), and in
1509 Chile for South American paleoecologists (Latorre et al., 2014). Other workshops train Data
1510 Stewards in how to prepare and upload data via Tilia. These Steward-oriented workshops have
1511 been provided as several-hour webinars, and as several-day workshops. We held Steward
1512 workshops at the American Quaternary Association (AmQua) 2016 biennial meeting and the
1513 European Pollen Database (2016). Instructional materials from these Steward- and user-
1514 oriented workshops and a template for developing new workshop materials are available on
1515 GitHub (<https://github.com/NeotomaDB/Workshops>). More workshops are planned.

1516 Neotoma collaborated with the University of Wisconsin Cartography Lab to hold a day-
1517 long Design Challenge mapping workshop in February 2016, bringing together students in
1518 cartographic design and paleoecologists to explore new approaches for visualizing paleodata
1519 (<https://www.earthcube.org/workspace/c4p/cartography-lab-design-challenge>). Neotoma helped
1520 organize a Community Development Workshop (also called a Paleodata Hackathon) in June
1521 2016, led by the Cyberinfrastructure for Paleogeosciences Research Coordination Network, in

1522 which interested scientific users and developers convened to explore, use, and link data from
1523 Neotoma, the Paleobiology Database, Macrostrat, and other resources.

1524 **Manuals and Workbooks.** The Neotoma Manual is available at <http://neotoma-manual.readthedocs.org/en/latest/>. A Tilia Manual is available at <http://tilia-manual.readthedocs.io/en/latest/>. Tilia help is also available from <https://www.tiliait.com/help/>.
1525 Sample workbooks are also available that provide users with standard code and fully realized
1526 workflows for common analyses of paleoecological data (<http://neotomadb.github.io>).

1527 **Educational Resources** are housed on the Neotoma website
1528 (https://www.neotomadb.org/education/category/higher_ed/) and at the Carleton SERC website
1529 (<http://serc.carleton.edu/neotoma/index.html>). These are varied, but most are interactive lab-
1530 based exercises, primarily geared towards college-level or high school courses. We welcome
1531 development of Neotoma-related educational materials and are happy to host or link to such
1532 content on the Neotoma home website. Interested contributors should contact the Neotoma
1533 Education and Outreach Committee.

1534 Several **Visualization Resources** are in development and are intended for a mixture of
1535 educational, outreach, and research purposes. The Flyover Country app (<http://fc.umn.edu>) is
1536 intended as an outreach and research tool that makes geological data available to any traveler
1537 with a smartphone (Loeffler et al., 2015), and showcases Neotoma fossil data holdings. An
1538 interactive map-based visualization, called Ice Age Mapper, is in development and is intended
1539 to replace the now-defunct Pollen Viewer (Williams et al., 2004). Pollen Viewer is no longer
1540 compliant with Java security standards and relied upon static images; Ice Age Mapper will
1541 dynamically link to Neotoma data holdings.

1542 **Issue Tracking and Code Repositories.** Most code for Neotoma is stored in GitHub, in
1543 the NeotomaDB organization (<https://github.com/NeotomaDB>). Individual code repositories
1544 within this organizational domain include the repository for minting DOIs
1545 (<https://github.com/NeotomaDB/AssignDOIs>), the *neotoma* R package
1546 (<https://github.com/NeotomaDB/neotoma>, forked from ropensci/neotoma), and code for mapping
1547 Neotoma to the DarwinCore standard (<https://github.com/NeotomaDB/DwC-Mapping>). The
1548 Neotoma API code is not currently open, but will be moved to an open repository in the near
1549 future (<https://github.com/NeotomaDB/Neotoma-API>). Issue tracking for Neotoma's
1550 Stratigraphic Diagrammer within Neotoma Explorer is available at
1551 <https://bitbucket.org/neotomadb/stratigraphicdiagram/issues> but we plan to migrate this to
1552 GitHub.
1553

Table S1: Summary of on-line resources

	URL	Description
Main Page	www.neotomadb.org	Main website for Neotoma Paleoecology Database
Data Access and Analysis	www.github.com/NeotomaDB	GitHub code repository for Neotoma DB and associated software
	http://apps.neotomadb.org/explorer/	Neotoma Explorer interface for finding, downloading, and quickly visualizing data
	http://api.neotomadb.org	APIs for programmatic access to Neotoma
	https://cran.r-project.org/web/packages/neotoma/index.html	CRAN repository and download source for <i>neotoma</i> R package
	https://github.com/ropensci/neotoma	GitHub code repository for <i>neotoma</i> R package
	http://www.neotomadb.org/snapshots	Archived snapshots of the Neotoma database
	https://dx.doi.org/10.6084/m9.figshare.3376393.v1	Example of an archived snapshot posted to figShare.
	https://data.neotomadb.org/datasets/1001/	Example of a digital object identifier (DOI) and associated landing page for an individual Neotoma dataset.
Data Contributions	https://www.neotomadb.org/data/category/contribution	Starting page for scientists interested in contributing data to Neotoma
	https://www.tiliait.com/	Website for downloading Tilia software and licenses.
Governance, Membership, Policy	http://www.neotomadb.org/about/category/governance	Neotoma's Governance Policy and By-Laws
	http://bit.ly/2tzjEsZ	List of Data Stewards (Google Docs)
	https://tinyurl.com/NeotomaMember	Webform to sign up as Neotoma member
	http://www.neotomadb.org/data/category/use	Data Use Policy

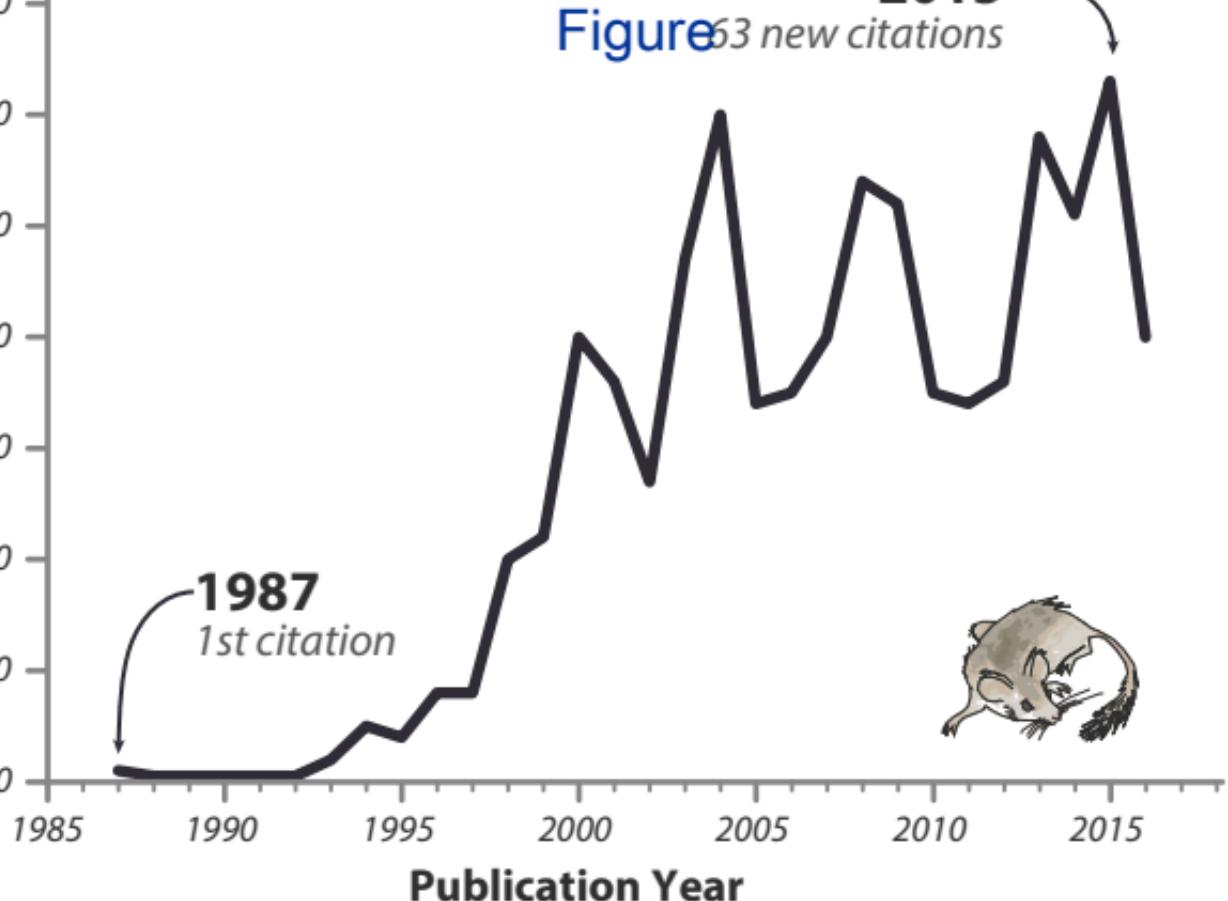
Help and Training Resources	<p>http://www.neotomadb.org/uploads/NeotomaManual.pdf</p> <p>http://neotoma-manual.readthedocs.org/en/latest/</p> <p>https://www.tiliait.com/help/</p> <p>http://tilia-manual.readthedocs.io/en/latest/</p> <p>https://github.com/NeotomaDB/FAU_NMAP_Import</p> <p>http://neotomadb.github.io</p>	<p>Neotoma Database Manual, with information about relational database</p> <p>A port of the Neotoma Database Manual to the readthedocs.org website, which</p> <p>Help page for Tilia software</p> <p>Port of the Tilia manual to readthedocs.org</p> <p>Example of a batch script for importing data into the Tilia .tlx format</p> <p>Sample scripts and analytical workflows making use of Neotoma data</p>
Teaching Resources	<p>https://www.neotomadb.org/education/category/higher_ed/</p> <p>http://serc.carleton.edu/neotoma/index.html</p>	<p>Access location for contributed teaching exercises and other educational researchers</p> <p>Educational resources for college and high school courses developed in partnership with the Science Education Research Center (SERC) at Carleton</p>
Allied and Linked Resources	<p>http://www.earthlifeconsortium.org/api_v1/ui/</p> <p>http://fc.umn.edu/</p> <p>http://www.gbif.org/</p> <p>https://globalpollenproject.org/</p> <p>(https://www.ncdc.noaa.gov/data-access/paleoclimatology-data</p>	<p>EarthLifeConsortium APIs that simultaneously search Neotoma and the Paleobiology Database.</p> <p>Flyover Country, a mobile app for discovering geological data during travel.</p> <p>Global Biodiversity Information Facility. We are in the process of exporting Neotoma data to GBIF.</p> <p>Global Pollen Project, a community platform for pollen images and identification</p> <p>NOAA Paleoclimatology Portal, searches will retrieve Neotoma data</p>

Figure 1

Click here to download
Figure 63 new citations



Number of New Papers Citing
Neotoma and Constituent DBs



Data Creators, Stewards

Biomarkers

Diatoms

Insects

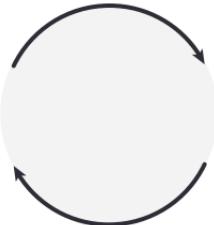
Ostracodes

Packrat
Middens

Pollen

Testate
Amoebae

Vertebrates

Contribute, curate
new data

Neotoma DB

New questions,
hypotheses, methodsProvide best practices &
recommended protocolsGenerate new questions,
hypotheses, methodsScientific drivers &
use cases

Paleoecologists

Archaeologists

Biogeographers

Ecologists

Educators

Paleoclimatologists

EarthCube

rOpenSci

DataOne

ICSU-WDS

ESIP

Informatics & Computer Scientists

Figure 3

[Click here to download Figure FigFX3_NeotomaGovernance_Rob_JWW.pdf](#)

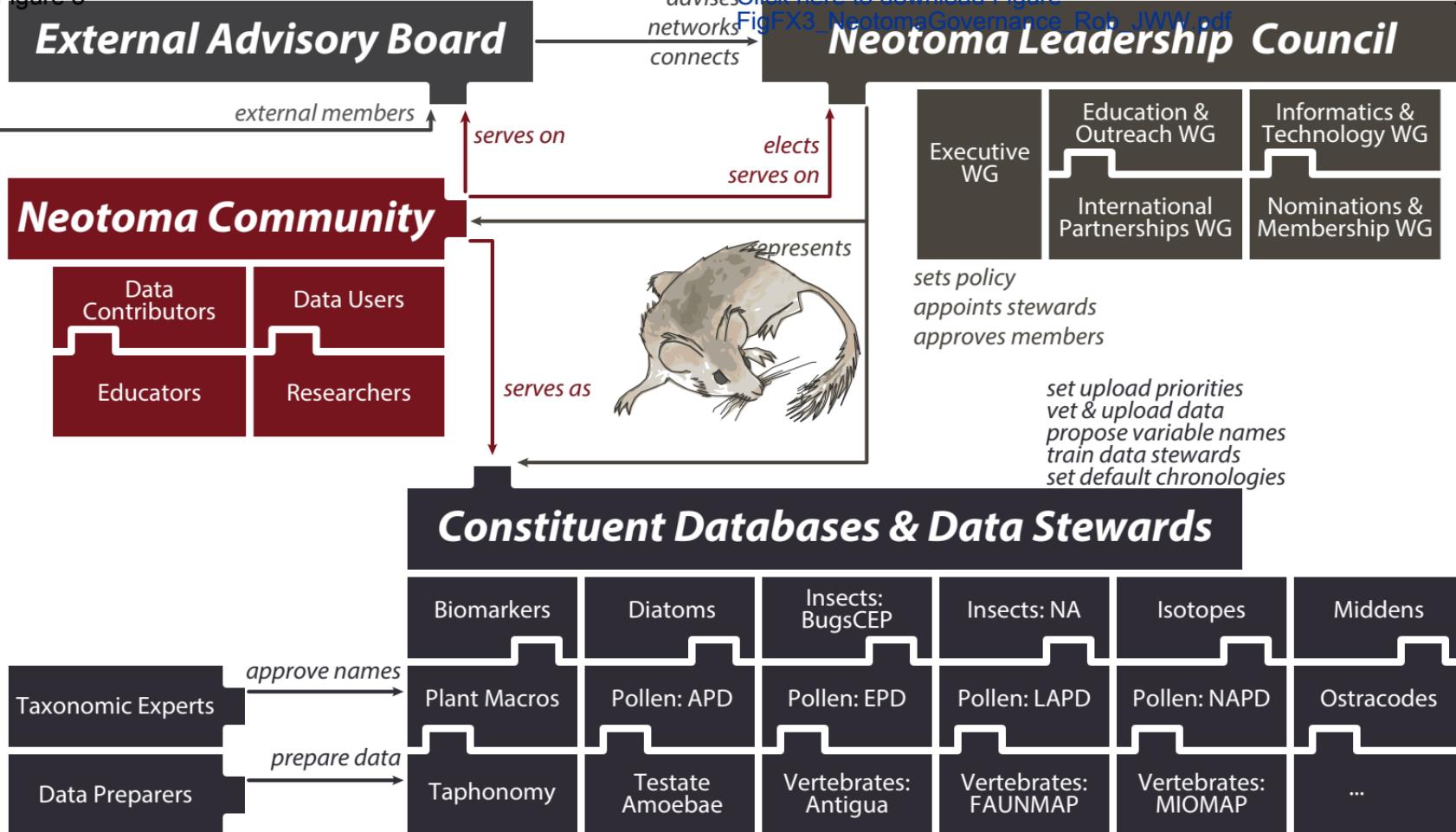


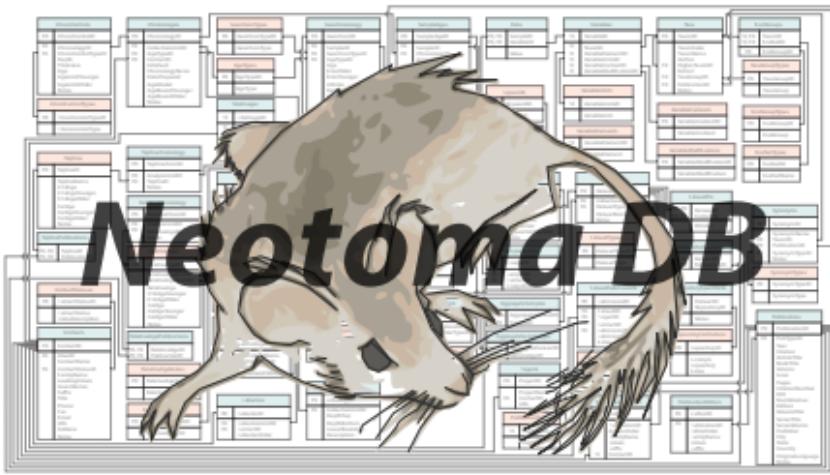
Figure 4

Click here to download Figure
FigFX4_NeotomaSoftwareEcosystem_Rob_JWW.pdf

Preparation, Validation, Submission, Revision



Archiving, Provisioning



Exploration, Discovery, Visualization, Retrieval



Figure 5

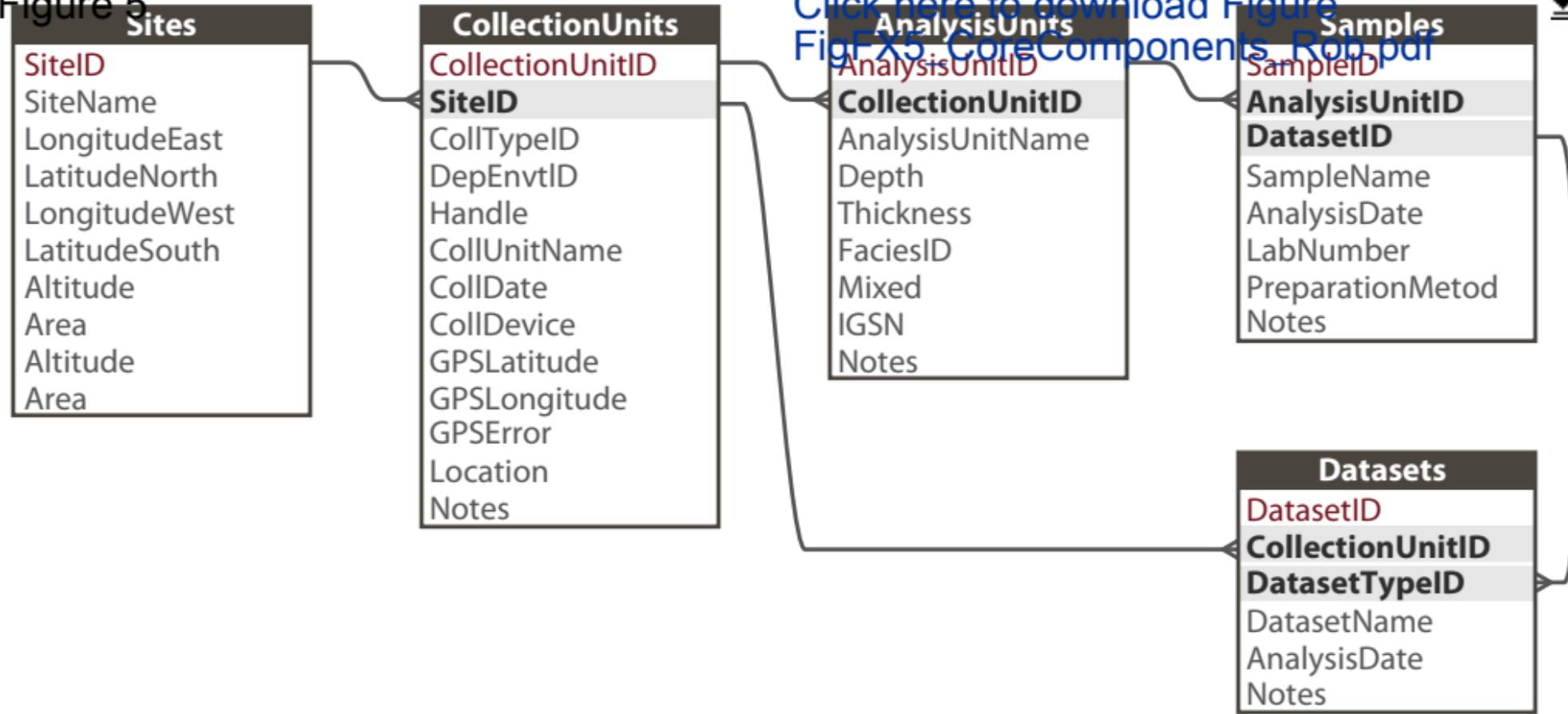
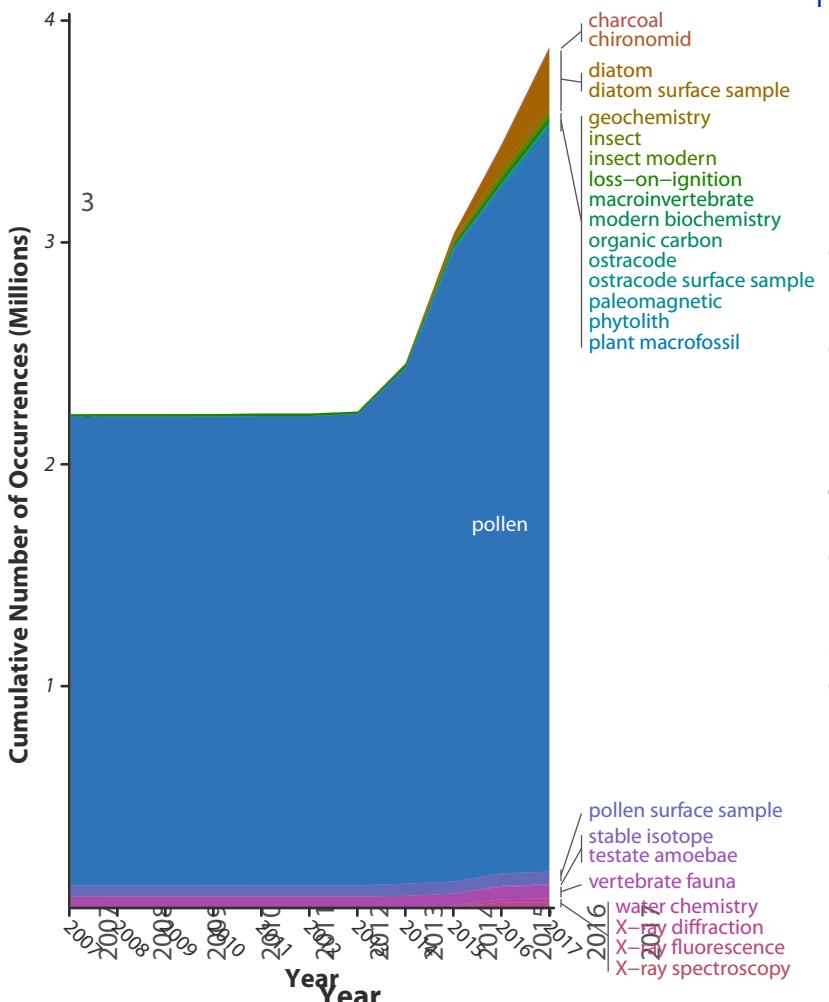


Figure 6



Click here to download Figure
FigFX6_CombinedGraphs_Rob.pdf

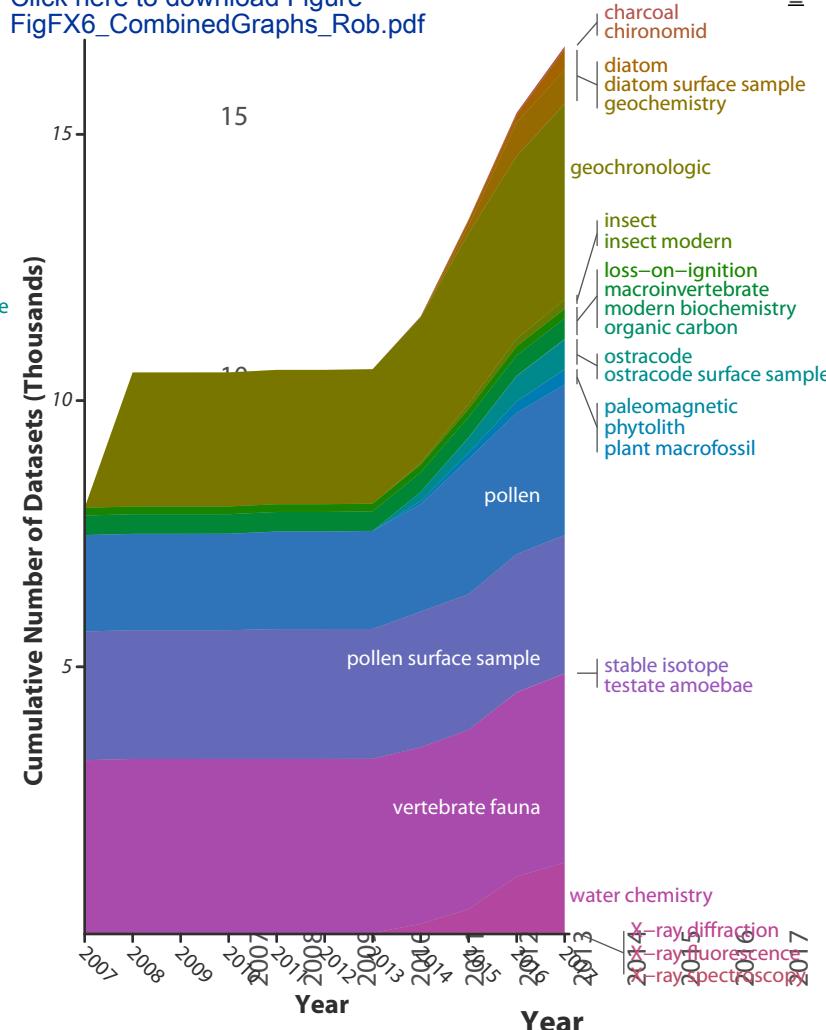


Figure 7
New Taxon[Click here to download Figure FigFX7_TiliaTaxonName.pdf](#)

Astrocsia

Select Taxa Group
Vascular plants

Insert New Taxon
 As Child
 As Sibling
Insert ↗

Insert Taxon
 As Child
 As Sibling
Insert ↗

Search: Phyllanthaceae

- > · Passifloraceae
- Phyllanthaceae
 - > · Phyllanthaceae subf. Antidesmatoideae
 - Phyllanthaceae subf. Phyllanthoideae
 - Amanoa
 - Andrachne
 - Astrocsia
 - Bridelia
 - Flueggea
 - > · Glochidion
 - Glochidion/Phyllanthus
 - Margaritaria
 - Margaritaria discoidea-type
 - > · Phyllanthus
 - > · Phyllanthus-type
 - Poranthera
 - Securinega
 - Securinega virosa-type
 - > · Picrodendraceae
 - > · Putranjivaceae
 - > · Rhizophoraceae
 - > · Salicaceae
 - > · Trigoniaceae

Locked Undo Close Validate

ID	Code	Name	Author	HigherID	Extinct	GroupID	PubID	Notes	EcolGroup	Validator	ValidDate
29927	Aoa	Astrocsia	B.L. Robinson & Millspaugh, 1905	9844		VPL	799		1:TRSH	44	2017-05...

Navigation icons: back, forward, search, etc.