1  **Digitization and the future of natural history collections**
2
3  Hedrick, B. P.*[a,b], Heberling, J. M.[+c], Meineke, E. K.[+b,d], Turner, K. G.[+e], Grassa, C. J.[d],
4  Park, D. S.[d], Kennedy, J.[d], Clarke, J. A.[f], Cook, J. A.[g], Blackburn, D.C. [h], Edwards, S.
5  V.[b], Davis, C. C.*[d]
6
7  [a]Department of Earth Sciences, University of Oxford, Oxford OX1 3AN United Kingdom
8  [b]Department of Organismal and Evolutionary Biology, Harvard University, Cambridge,
9  MA 02138
10  [c]Section of Botany, Carnegie Museum of Natural History, Pittsburgh, PA 15213
11  [d]Harvard University Herbaria, Harvard University, Cambridge, MA 02138
12  [e]Department of Biology, Pennsylvania State University, University Park, PA 16802
13  [f]Jackson School of Geosciences, University of Texas Austin, Austin, TX 78712
14  [g]Department of Biology, University of New Mexico, Albuquerque, NM 87131
15  [h]Florida Museum of Natural History, University of Florida, Gainesville, FL 32611
16
17  *Corresponding authors: bphedrick1@gmail.com, cdavis@oeb.harvard.edu
18  [+]Contributed equally
19
20  Keywords: digitization, herbaria, natural history collections, specimens, Anthropocene,
21  baselines
22

## Abstract

Natural history collections (NHCs) are the foundation of historical baselines for assessing

anthropogenic impacts on biodiversity. Along these lines, the online mobilization of

specimens via digitization–the conversion of specimen data into accessible digital

content–has greatly expanded the use of NHC collections across a diversity of

disciplines. We broaden the current vision of digitization (Digitization 1.0)–whereby

specimens are digitized within NHCs–to include new approaches that rely on digitized

products rather than the physical specimen (Digitization 2.0). Digitization 2.0 builds upon

the data, workflows, and infrastructure produced by Digitization 1.0 to create digital-only

workflows that facilitate digitization, curation, and data linkages, thus returning value to

physical specimens by creating new layers of annotation, empowering a global

community, and developing automated approaches to advance biodiversity discovery and

conservation. These efforts will transform large-scale biodiversity assessments to address

fundamental questions including those pertaining to critical modern issues of global

change.

38    **I. The relevance and importance of digitization**

39    Anthropogenic impacts, including urbanization, globalization, and climate change, are

40    rapidly transforming our world. Despite our best efforts, however, quantifying the biotic

41    impacts of human activity has been challenging, as evidenced by the difficulty of

42    delimiting the onset of the Anthropocene (Lewis and Maslin 2015). Part of this

43    uncertainty stems from a lack of historical data that track biotic change over time.

44    However, natural history collections (NHCs), with their broad taxonomic, geographic,

45    and temporal scope, offer a key solution to this impasse. In the past twenty years, there

46    has been a dramatic increase in the use of NHCs for assessing a wide variety of scientific

47    questions (Suarez and Tsutsui 2004, Pyke and Ehrlich 2010, Park and Potter 2015,

48    Meineke et al. 2018, 2019). Indeed, they have emerged as one of the best resources for

49    establishing biological baselines to understand the impacts of, for example, the origins of

50    agriculture, the industrial revolution, the development of nuclear armaments, and more

51    generally the influence and acceleration of anthropogenic change on biodiversity (Moritz

52    et al. 2008, Johnson et al. 2011, Lister 2011, Funk 2018, Nelson and Ellis 2018).

53

54    Most large NHCs provide specimen data to researchers and the public by mobilizing

55    searchable collection databases online. We assert that these mobilized collections are

56    among the most important advances in museum curation in the past century, significantly

57    opening access to NHCs and greatly stimulating large-scale analyses that span novel

58    academic and societal enterprises. These resources are connecting diverse scholarly

59    domains, propelling a new generation of scientists forward, and removing financial,

60    sociological, institutional, and academic obstacles preventing access to these materials

61    (Drew et al. 2017, Sweeney et al. 2018). In short, digitizing a specimen–translating

62    metadata associated with a physical specimen object into flexible digital data formats–

63    increases the value of the physical specimen exponentially.

64

65    Here, we present an ambitious, two-pronged vision for digitization, which we term

66    Digitization 1.0 and Digitization 2.0. Digitization 1.0 represents the ongoing push to

67    create digital images and related content directly from physical voucher specimens;

68    Digitization 2.0, in contrast, relates exclusively to data gathering, tasks, or workflows

69    derived from digitized products of Digitization 1.0 rather than from the physical

70    specimens themselves (figure 1). In addition to the vast expansion and online aggregation

71    of these mobilized collections to create a truly global digital NHC, Digitization 2.0 offers

72    the promise of also shifting and growing the workforce and public who interface with

73    these objects to accelerate the progress of digitization.

74

75    **II. Digitization 1.0: The Past, Present, and Future**

76    Digitization of NHCs began with the overarching goal of documenting specimen

77    inventory and facilitating research by transcribing label information into centralized,

78    searchable databases as described recently by Nelson and Ellis (2018). These efforts have

79    given rise to Digitization 1.0, which has been widely embraced and continues to be

80    infused with innovation. Digital representations generated through Digitization 1.0

81    include specimen images and direct transcriptions of specimen metadata from

82    handwritten or printed collection catalogs or labels, including for example details on

83    coloration or measurements. As part of this effort, NHCs have generated millions of

84    digital representations of physical vouchers and have devised numerous technological

85    innovations to facilitate efficient data generation, including conveyor belt and robotic

86    imaging techniques for mass digitization of specimens (Tegelberg et al. 2014, Sweeney et

87    al. 2018). More recent next generation technologies, including photogrammetry, laser-

88    scanning, and computed tomography, create far richer digital representations of

89    specimens than can be visualized by eye or with standard microscopy (figure 2). Given

90    that large portions of most NHCs still remain unavailable in digital format, the

91    innovations and efforts within Digitization 1.0 will continue well into the future, likely

92    for decades. In the subsections below, we outline Digitization 1.0 through the lens of

93    digitization workflows, strategic prioritization, and solutions to impediments.

94

95    ***Digitization workflows and linking data***–The practice of digitization is broadly

96    consistent among projects and organismal groups, in so much as each specimen is

97    represented by textual metadata from labels or catalogs and typically digital two-

98    dimensional images, but increasingly also three-dimensional representations and audio or

99    video recordings where relevant. There exists great variation in specimen size, storage

100   conditions (e.g., fluid-preserved, microscope slides, dry storage), dimensionality (2D

101   versus 3D representation), and detail associated with specimens, not to mention widely

102   varying practices in specimen collection and curation across taxonomic domains and

103   institutions. This heterogeneity of collections and institutional policies and priorities thus

104   creates challenges to efficient mass imaging and gathering of metadata. However, at

105   minimum, digitization workflows should attempt to integrate all available specimen

106   metadata into digitization efforts and appropriately link these data to their associated

107    physical voucher specimens. Beyond traditional linkages, non-traditional metadata

108    associated with the specimen include biotic (e.g., mass) and abiotic data (e.g., climate),

109    media (e.g., video and audio recordings), community- and population-level metadata

110    (e.g., abundance), species observations in the field, and genetic samples (i.e., the

111    "extended specimen" *sensu* Webster 2017). Much of these digital data are served in part

112    or in their entirety via online collection databases (e.g., Arctos, Specify, Symbiota, EMu)

113    or in data aggregators (e.g., iDigBio, Global Biodiversity Information Facility–GBIF,

114    Botanical Information and Ecology Network–BIEN). Linking voucher specimens to these

115    new data layers generated post collection is important and has been facilitated by

116    associating URLs, data accession numbers, digital object identifiers (DOI), or ARKs with

117    specimen records in collection databases. In addition, trait data can be incorporated into

118    specimen records using extensions to the Darwin Core Archives (Yost et al. 2018). For

119    the next generation of collections, protocols are under development to expand the

120    digitization workflow to the collecting event itself (Heberling and Issac 2018).

121

122    ***Developing digitization priorities***–Given the limited resources available to many NHCs,

123    it is necessary to establish priorities for specimen digitization. Specimens at risk of

124    degradation, such as rare or fragile fossils, and those representing rare or threatened

125    species and habitats are candidates for high priority digitization. Further, efforts should

126    focus on specimens with rich associated metadata from the collection event. A growing

127    number of species are imperiled, and conservation biologists are increasingly reliant on

128    NHCs for baseline data to understand species ranges and climatic tolerances for assessing

129    future changes (Lister 2011). Distributing information for these rare or threatened taxa to

130    conservation biologists is increasingly critical to these species' management and survival

131    (MacDougall et al. 1998, Nualart et al. 2017). Finally, taxa representing a breadth of

132    evolutionary history or unique adaptations are important for research on phenotypic

133    evolution, community ecology, and biologically inspired design. We suggest that such

134    specimens have high priority for digitization.

135

136    Owing to the varying effort required by different digitization strategies (e.g., label data,

137    images, 3D reconstructions), data types that serve the largest diversity of use cases should

138    also be prioritized. For instance, key information including taxon name, collection

139    locality, and date can be captured relatively efficiently and can facilitate assessments of

140    species distributions through time. Rapidly expanding areas of research including

141    phenology (e.g., Primack et al. 2004, Willis et al. 2017), large-scale taxonomic

142    inventories (e.g., Cardoso et al. 2017), and morphometric investigations (e.g., Hedrick et

143    al. 2015), rely on such label data and data from post-digitization enhancement (Sweeney

144    et al. 2018). For example, in one of the first studies to demonstrate how historic

145    specimens can be used to quantify the biotic effects of climate change, Primack et al.

146    (2004) used flowering plant specimens collected between 1885 and 2003 in the greater

147    Boston (USA) area to demonstrate that plants were flowering up to eight days earlier in

148    recent years than in the early years of the 20[th] century. The utility of such diverse data

149    (e.g., geographic location, flowering date, anatomical measurements) is important to a

150    wide array of researchers and should be prioritized. Additionally, we feel it is best to only

151    apply more complex, holistic digitization methods on a key subset of data-rich specimens

152    as has been recently demonstrated in the openVertebrate (oVert) Thematic Collection

7

153    Network (Blackburn et al., NSF Abstract #1701714). Increasing the magnitude of the

154    collection of media files (e.g., photogrammetry of bird skins, nuts, etc.) for this subset of

155    data via new pipelines and technological advances will be critical to this effort.

156

157    ***Past impediments and future solutions*–**Despite the success of Digitization 1.0, this

158    initiative has identified three issues that must be addressed to maximize efficiency of

159    information retention and distribution. First, museums are obligated to manage, store, and

160    steward additional digital data associated with their physical collections. However, the

161    act of digitization entails significant challenges since it requires sustainably curating both

162    the physical objects and rapidly emerging digital datasets. This issue will necessitate the

163    development of new tools, that centralized aggregators assume increasing responsibility,

164    and will require increased funding in the near future (see Digitization 2.0 below).

165

166    Second, there is concern that large aggregators aimed at connecting researchers with

167    NHCs (e.g., GBIF, iDigBio) (Edwards 2004) remove NHCs from the attribution chain.

168    NHCs are frequently funded on their research relevance, which is determined both from

169    within and outside institutions. When researchers view specimen images or harvest

170    metadata from aggregators, NHCs that contribute these data often receive little to no

171    credit (Rouhan et al., 2017). A mechanism for referencing these source collections needs

172    to be embedded in the publication process that requires that NHCs be acknowledged and

173    notified when publications incorporate their data. A viable solution to this problem is to

174    mint a digital object identifier (DOI) for a digitized specimen and establish a reporting

175    mechanism for collections to be alerted when their specimens have been cited.

176    Automating this attribution pipeline as part of the digitization workflow better ensures

177    that NHCs receive credit for stewarding both voucher specimens and also digitized data,

178    which is likely to stimulate NHCs to embrace open-access policies for their data.

179

180    Third, digitized data are inconsistently and redundantly spread across multiple databases

181    at different scales. NHCs often have their own databases, but some data are additionally

182    deposited in regional databases, taxon-specific databases, and national and international

183    data aggregators. This data dispersion causes information to be input/archived

184    redundantly such that each database has a variant of the post-digitization metadata,

185    leading aggregators to archive either inconsistent or duplicated copies of the same

186    primary data. This problem can be partially circumvented by more communication

187    among data aggregators, as well as between NHCs and aggregators. Algorithms linking

188    specimen numbers between aggregators could ensure that post-digitization enhancement

189    metadata are transferred to all aggregators mentioning particular specimens by unique

190    identifiers such as the specimen-based occurrenceID. This is done internally at iDigBio

191    through the iDigBio Record API, which retains current and previous iterations of a

192    specimen's data.

193

194    **III. Digitization 2.0: charting a road map for the future**

195    Unlike Digitization 1.0, which directly utilizes the physical specimen, Digitization 2.0

196    instead utilizes the digitized product from Digitization 1.0 for generating additional data

197    and metadata (figure 1). Digitization 2.0 is powered by the online aggregation of these

198    resources and enables digitization to assume new forms and engage vast new workforces.

199    As we outline below, Digitization 2.0 is already well underway and holds tremendous

200    promise. It includes semi- or fully automated data recording from digitized specimens,

201    which stimulates research and returns value to the physical specimen. Additionally,

202    Digitization 2.0 entails a shift in the workforce engaged in collections science and

203    stewardship. Finally, Digitization 2.0 leverages NHC resources to create trait databases,

204    either from aggregating and better indexing existing metadata or by allowing researchers

205    or citizen scientists to associate trait annotations with images served from NHC

206    databases.

207

208    ***Innovative tools for automating digitization: machine learning and neural networks***–

209    Given the massive number of specimen images in digital databases with minimal data, an

210    important first step is to better automate data transcription to augment these skeletal

211    records. The enormity of this task is quickly becoming impossibly large for collections

212    staff to manage without automation, especially considering that funding for NHCs has

213    been decreasing (Thiers 2018). In recent years, machine learning applications utilizing

214    convolutional neural networks have achieved stunning levels of performance in computer

215    vision tasks including image detection and classification (Sudholt and Fink 2016). Neural

216    networks have previously demonstrated promising results for handwriting recognition

217    systems, which could easily be applied to automated label transcription. These forms of

218    innovative technology, which have been applied to medical diagnoses, speech

219    recognition, and driverless cars, are now permeating NHCs (Schuettpelz et al. 2017) and

220    are likely to be enormously useful when trained on existing databases of handwriting

221    samples (Krishnan et al. 2016), as well as those from transcribed labels generated through

222   Digitization 1.0. These models can be further trained using existing semantic field

223   constraints to much more effectively parse specimen metadata into appropriate database

224   fields. Beyond capturing essential minimal data records in an automated manner, neural

225   networks have recently been implemented to accomplish far more sophisticated tasks

226   than text transcription (Wilf et al. 2016, Schuettpelz et al. 2017, Funk 2018). Wilf et al.

227   (2016), for example, used computer vision to classify fossil leaf images based on leaf

228   shape and venation with high accuracy. This proved not only to be an efficient protocol

229   for classifying images, but also discovered previously unidentified morphological

230   landmarks potentially useful for species identification and for understanding important

231   evolutionary and ecological innovations not previously documented. The community is

232   now ready for deeper exploration of minimal metadata capture using semi- to total

233   automation.

234

235   Further, the declining number of taxonomists in the global workforce severely impacts

236   our ability to address key questions concerning biodiversity in the face of global change

237   (Hopkins and Freckleton 2002). Combining taxonomists' expertise (past and present)

238   with student and public training and increased automation will facilitate enhanced

239   specimen curation, and greatly enable biodiversity discovery. Continued robust support

240   for taxonomic research and training is essential. However, given the enormity of the task

241   at hand, and the limited time for this effort, we believe that addressing many taxonomic

242   problems of identification, particularly for well-known groups of organisms, could be

243   greatly facilitated by automation, such as has been demonstrated through Kurator (Dou et

244   al. 2012). Reasonably successful early efforts are underway to machine-learn and

245    automatically identify large sub-collections of insects (e.g., butterflies) (Schermer and

246    Hogeweg 2018). Although simple taxonomic identification may seem rudimentary, it is

247    the foundation of all biological research, and in many groups remains problematic. For

248    example, it is estimated that more than 50% of tropical plant specimens in NHCs are

249    incorrectly identified (Goodwin et al. 2015). Together with the training of more expert

250    taxonomists and organismal biologists, the widespread use of neural networks to identify

251    specimens and target groups that need attention would enhance collection utility for

252    research, teaching, and management and further motivate the discovery and description

253    of new species.

254

255    ***Expansion of the digitization workforce***–Expanding digitization to involve a global

256    workforce is now possible and is a major advancement of Digitization 2.0 and is

257    motivated by the increasingly global accessibility of NHCs. These new workforces can

258    be developed to supplement existing NHC staff, especially to include enhanced

259    digitization from the millions of images in databases that have limited associated

260    metadata. One obvious group to engage in this effort are citizen scientists. NHCs

261    associated with museums typically have departments devoted to public outreach, which

262    can easily be tapped for aid, helping collections staff with the task of digitization while

263    simultaneously providing the public with ownership and agency. Using citizen science in

264    this manner has been fruitful in numerous contexts including the transcriptions of label

265    data, georeferencing, and physical specimen annotations (Hill et al. 2012, Ballard et al.

266    2017, Ellwood et al. 2015, 2017). For example, *CrowdCurio–Thoreau's Field Notes*, an

267    online crowdsourcing platform has successfully facilitated climate change studies from

268    thousands of herbarium specimens utilizing thousands of non-expert crowdsourcers

269    (Willis et al. 2017). Quality control is always a concern in large-scale citizen science

270    projects (Willis et al. 2017, Zhou et al. 2018) and thus an easy-to-use graphical user

271    interface clearly demonstrating to the public how and what to digitize will be necessary

272    (e.g., Notes for Nature), as has been accomplished in several research-based projects

273    (Chang and Alfaro 2016, Cooney et al. 2017, Willis et al. 2017). Increasingly, such

274    citizen science efforts are being supplemented by machine-based learning as well (Unger

275    et al. 2016, Wilf et al. 2016, Schuettpelz et al. 2017). For instance, crowdsourced data can

276    potentially provide reliable and rapid data for training and testing machine learning

277    models, creating a positive feedback loop propelling digitization forward.

278

279    ***Layers of trait annotations***–Traits of organisms are fundamental for documenting

280    biodiversity but also for understanding how organisms evolve and respond to changing

281    environments. Building on investments in creating digital NHCs, there is now increasing

282    demand for creating and associating new layers of trait data to these collections. For

283    some taxa, these biological data are already captured in the digitized text of a specimen

284    record (e.g., Darwin Core fields: "organismRemarks"). In mammals and birds, it is

285    common to have measurements on the mass and length of both the whole specimen and

286    parts of the specimen (e.g., testes length, wing length). The aggregation of traits from

287    both the initial collecting event and new annotations will stimulate a wealth of questions

288    and generate a better understanding of global biodiversity through the development of

289    standardized trait vocabularies (Kissling et al. 2018). For example, recently developed

290    data-processing tools for the data aggregator VertNet standardized more than 1.5 million

291    measurements for vertebrates using digital data from collections (Guralnick et al. 2016).

292    Users can now search those specimen records by mass and length, as well as download

293    harmonized trait data associated with individual specimens. The latter allows for new

294    explorations of trait variation within and across species, including spatial and temporal

295    patterns in traits associated with specimens that have collecting dates and georeferenced

296    localities (Riemer et a. 2018). By expanding this framework to annotate traits to

297    specimens and utilizing online platforms for even 3D representations of specimens,

298    NHCs can facilitate the capture of not only simple traits, ranging from specimen length to

299    the presence of a flower, but also more complex traits requiring more sophisticated

300    representation (e.g., virtual automated dissection of the vertebrate nervous system).

301

302    **IV. Concluding thoughts**

303    Digitization facilitates the democratizing of collections-based research and is essential to

304    establishing and evaluating biological baselines to assess the impacts of climate change,

305    land use changes, species invasions, and the current mass extinction. It allows for the

306    mining of specimen data in much the same way that we explore organismal genomes.

307    The key to further developing Digitization 1.0 and establishing Digitization 2.0 lies in

308    building upon what the research, funding, and policy communities have learned in the

309    several decades since the initiation of this endeavor. Data-rich NHC specimens are useful

310    and provide unique perspectives on the diversity and distribution of a given taxon.

311    However, if a specimen is not searchable, it will likely not be found or studied despite its

312    potential use. We are already witnessing the fruits of the synergy between Digitization

313    1.0 and 2.0. This will no doubt expand dramatically in the coming decades to involve

314    new domains, new questions, and new audiences that are not yet realized (or even

315    imagined). Only with creativity and improved techniques, including automated and semi-

316    automated methods, a better distributed digitization workload making use of new

317    technologies and workforces, and conscientious attention to the attribution chain, will

318    researchers be best able to track ongoing biodiversity change from all existing data.

319    Moreover, even as new technologies and digitization techniques emerge, we will need to

320    always return to physical specimens, in ways that are unimaginable now, to generate

321    novel data to better understand our changing planet. Although we stress the importance of

322    improved methods and practices for digitization, the active collection and continued

323    curation of physical specimens by expert biologists remains the central pillar supporting

324    advancements in evolutionary biology and conservation represented so importantly by

325    NHCs.

326

327    **Box 1:** *Estimating the size and scale of a global digitization effort*– Digitization 1.0 has

328    resulted in the mobilization of millions of specimen records and has created the

329    momentum for a massive, global digitization effort. To better establish target goals and

330    evaluate the success of this effort (e.g., estimating the proportion of specimen records that

331    have been digitized and mobilized online), obtaining accurate estimates of the number of

332    specimens housed in NHCs is necessary. Extrapolations from digitized content indicate

333    that roughly 2.5–3 billion specimens are housed in NHCs worldwide (O'Connell et al.

334    2004, Krishnan et al. 2016). However, more robust assessments of global specimen

335    numbers, including geographic and taxonomic distribution, are required to facilitate

336    thoughtful assessments of collection bias to better target digitization priorities (Meyer et

337   al. 2016). Making robust size estimates are particularly relevant as vended solutions are

338   utilized to achieve digitization milestones, including the mobilization of entire collections

339   like those at the Muséum National D'Histoire Naturelle (France), Naturalis

340   (Netherlands), and the Smithsonian Institution (US) (Rogers 2016, Le Bras et al. 2017).

341   Along these lines, a test case example to illustrate such an effort on a smaller scale comes

342   from the Harvard University Herbaria (HUH), which has been thought to contain 5.5

343   million specimens. Targeted subsampling of the HUH vascular plant collection facilitated

344   accurate estimates (with confidence intervals) of total specimen collection numbers and

345   their geographic distribution (figure 3A). Once the total number of specimens in NHCs

346   have been accurately quantified, it is necessary to establish the percentage of specimen

347   collection records that have been digitally mobilized.

348       Because imaging and serving metadata-rich collection information online requires

349   a large financial investment, as well as human labor, its impacts on research should be

350   documented and acknowledged. The most powerful outcomes of digitization would be

351   better characterized by relating these various forms of data usage to one another to

352   explore how digitization increases specimen usage. Along these lines, data relevant to

353   describing the scientific impact of physical specimens (pre-digitization), such as loans

354   and museum visits, remain largely confined to physical collection logbooks, thus limiting

355   assessment of the impact of Digitization 1.0 (figure 3B). Such efforts would allow us to

356   begin to understand the ways that digitization stimulates increased visitation and use of

357   the actual physical versus digital collection (figure 3C). As a community, we must be

358   better prepared to track and assess these questions.

359

## Acknowledgements

## References Cited

Ballard HL, Robinson LD, Young AN, Pauly GB, Higgins LM, Johnson RF, Tweddle JC.
2017. Contributions to conservation outcomes by natural history museum-led
citizen science: Examining evidence and next steps. *Biological Conservation* 208:
87–97.

Cardoso D, Särkinen T, Alexander S, Amorim AM, Bittrich V, Celis M, Daly DC,
Fiaschi P, Funk VA, Giacomin LL, Goldenberg R. 2017. Amazon plant diversity
revealed by a taxonomically verified species list. *PNAS* 114: 10695–10700.

Chang J, Alfaro ME. 2016. Crowdsourced geometric morphometrics enable rapid large-
scale collection and analysis of phenotypic data. *Methods in Ecology and
Evolution* 7: 472–482.

Comoglio F, Fracchia L, Rinaldi M. 2013. Bayesian inference from count data using
discrete uniform priors. *PLoS One* 8: e74388.

Cooney CR, Bright JA, Capp EJ, Chira AM, Hughes EC, Moody CJ, Nouri LO, Varley
ZK, Thomas GH. 2017. Mega-evolutionary dynamics of the adaptive radiation of
birds. *Nature* 542: 344–347.

Dou L, Cao G, Morris PJ, Morris RA, Ludäscher B, Macklin JA, Hanken J. 2012.
Kurator: A Kepler package for data curation workflows. *Procedia Computer
Science* 9: 1614–1619.

Drew JA, Moreau CS, Stiassny MLJ. 2017. Digitization of museum collections holds the
potential to enhance researcher diversity. *Nature Ecology & Evolution* 1: 1789.

17

Edwards JL. 2004. Research and societal benefits of the global biodiversity information facility. *Bioscience* 54: 485–486.

Ellwood ER, Dunckel BA, Flemons P, Guralnick R, Nelson G, Newman G, Newman S, Paul D, Riccardi G, Rios N, Seltmann KC. 2015. Accelerating the digitization of biodiversity research specimens through online public participation. *Bioscience* 4: 383–396.

Ellwood ER, Crimmins TM, Miller-Rushing AJ. 2017. Citizen science and conservation: Recommendations for a rapidly moving field. *Biological Conservation* 208: 1–4.

Funk VA. 2018. Collections-based science in the 21st century. *Journal of Systematics and Evolution* 56: 175–193.

Goodwin ZA, Harris DJ, Filer D, Wood JR, Scotland RW. 2015. Widespread mistaken identity in tropical plant collections. *Current Biology* 25: 1066–1067.

Guralnick RP, Zermoglio PF, Wieczorek J, LaFrance R, Bloom D, Russell L. 2016. The importance of digitized biocollections as a source of trait data and a new VertNet resource. *Database* Volume 2016.

Heberling JM, Isaac BL. 2018. iNaturalist as a tool to expand the research value of museum specimens. *Applications in Plant Sciences* 6: e1193.

Hedrick BP, Manning PL, Lynch ER, Cordero SA, Dodson P. 2015. The geometry of taking flight: Limb morphometrics in Cretaceous theropods. *Journal of Morphology* 276: 152–166.

Hill A, Guralnick R, Smith A, Sallans A, Gillespie R, Denslow M, Gross J, Murrell Z, Conyers T, Oboyski P, Ball J. 2012. The notes from nature tool for unlocking biodiversity records from museum records through citizen science. *ZooKeys* 209: 219.

Hopkins GW, Freckleton RP. 2002. Declines in the numbers of amateur and professional taxonomists: Implications for conservation. *Animal Conservation* 5: 245–249.

Johnson KG, Brooks SJ, Fenberg PB, Glover AG, James KE, Lister AM, Michel E, Spencer M, Todd JA, Valsami-Jones E, Young JR. 2011. Climate change and biosphere response: Unlocking the collections vault. *Bioscience* 61: 148–153.

Kissling WD, Walls R, Bowser A, Jones MO, Kattge J, Agosti D, Amengual J, Basset A, Van Bodegom PM, Cornelissen JH, Denny EG. 2018. Towards global data products of essential biodiversity variables on species traits. *Nature Ecology & Evolution* 2: 1531–1540.

Krishnan P, Dutta K, Jawahar CV. 2016. Deep feature embedding for accurate recognition and retrieval of handwritten text. *Frontiers in Handwriting Recognition (ICFHR)* 15th International Conference: 289–294.

Le Bras G, Pignal M, Jeanson ML, Muller S, Aupic C, Carré B, Flament G, Gaudeul M, Gonçalves C, Invernón VR, Jabbour F. 2017. The French Muséum National D'histoire Naturelle vascular plant herbarium collection dataset. *Scientific Data* 4: 170016.

Lewis SL, Maslin MA. 2015. Defining the Anthropocene. *Nature* 519: 171.

Lister AM. 2011. Natural history collections as sources of long-term datasets. *Trends in Ecology & Evolution* 26: 153–154.

MacDougall AS, Loob JA, Claydenc SR, Goltzd JG, Hindse HR. 1998. Defining conservation priorities for plant taxa in southeastern New Brunswick, Canada using herbarium records. *Biological Conservation* 86: 325–338.

440    Meineke EK, Davis CC, Davies TJ. 2018. The unrealized potential of herbaria for global
441         change biology. *Ecological Monographs* 88: 505–525.
442    Meineke EK, Davies TJ, Daru BH, Davis CC. 2019. Biological collections for
443         understanding biodiversity in the Anthropocene. *Philosophical Transactions of*
444         *the Royal Society B: Biological Sciences* 374: 20170386.
445    Meyer C, Weigelt P, Kreft H. 2016. Multidimensional biases, gaps and uncertainties in
446         global plant occurrence information. *Ecology Letters* 19: 992–1006.
447    Moritz C, Patton JL, Conroy CJ, Parra JL, White GC, Beissinger SR. 2008. Impact of a
448         century of climate change on small-mammal communities in Yosemite National
449         Park, USA. *Science* 322: 261–264.
450    Nelson G, Ellis S. 2018. The history and impact of digitization and digital data
451         mobilization on biodiversity research. *Philosophical Transactions of the Royal*
452         *Society B*. 374: 20170391.
453    Nualart N, Ibáñez N, Soriano I, López-Pujol J. 2017. Assessing the relevance of
454         herbarium collections as tools for conservation biology. *The Botanical Review* 83:
455         303–325.
456    O'Connell AFJ, Gilbert AT, Hatfield JS. 2004. Contribution of natural history collection
457         data to biodiversity assessment in national parks. *Conservation Biology* 18: 1254–
458         1261.
459    Park DS, Potter D. 2015. Why close relatives make bad neighbours: phylogenetic
460         conservatism in niche preferences and dispersal disproves Darwin's naturalization
461         hypothesis in the thistle tribe. *Molecular Ecology* 24: 3181–3193.
462    Primack D, Imbres C, Primack RB, Miller-Rushing AJ, del Tredici P. 2004. Herbarium
463         specimens demonstrate earlier flowering times in response to warming in Boston.
464         *American Journal of Botany* 91: 1260–1264.
465    Pyke GH, Ehrlich PR. 2010. Biological collections and ecological/environmental
466         research: A review, some observations and a look to the future. *Biological*
467         *Reviews* 85: 247–266.
468    Riemer K, Guralnick RP, White EP. 2018. No general relationship between mass and
469         temperature in endothermic species. *eLife 7*: e27166.
470    Rogers N. 2016. Museum drawers go digital. *Science* 352: 762–765.
471    Rouhan G, Dorr LJ, Gautier L, Clerc P, Muller S, Gaudeul M. 2017. The time has come
472         for natural history collections to claim co-authorship of research articles.
473         *Taxon* 66: 101–1016.
474    Schermer M, Hogeweg L. 2018. Supporting citizen scientists with automatic species
475         identification using deep learning image recognition models. *Biodiversity*
476         *Information Science and Standards* 2: e25268.
477    Schuettpelz E, Frandsen PB, Dikow RB, Brown A, Orli S, Peters M, Metallo A, Funk
478         VA, Dorr LJ. 2017. Applications of deep convolutional neural networks to
479         digitized natural history collections. *Biodiversity Data Journal* 5: e21139.
480    Suarez AV, Tsutsui ND. 2004. The value of museum collections for research collections.
481         *Bioscience* 54: 66–74.
482    Sudholt S, Fink GA. 2016. PHOCNet: A deep convolutional neural network for word
483         spotting in handwritten documents. *Frontiers in Handwriting Recognition*
484         *(ICFHR)* 15th International Conference: 277–282.

485  Sweeney PW, Starly B, Morris PJ, Xu Y, Jones A, Radhakrishnan S, Grassa CJ, Davis
486      CC. 2018. Large-scale digitization of herbarium specimens: Development and
487      usage of an automated, high-throughput conveyor system. *Taxon* 67: 165–178.
488  Tegelberg R, Mononen T, Saarenmaa H. 2014. High-performance digitization of natural
489      history collections: Automated imaging lines for herbarium and insect specimens.
490      *Taxon* 63: 1307–1313.
491  Thiers B. 2018. Using data from index herbarium to assess threats to the world's
492      herbaria. *Biodiversity Information Science and Standards* 2: e26440.
493  Unger J, Merhof D, Renner S. 2016. Computer vision applied to herbarium specimens of
494      German trees: Testing the future utility of the millions of herbarium specimen
495      images for automated identification. *BMC Evolutionary Biology* 16: 248.
496  Webster M. 2017. The extended specimen. In M. Webster [ed.], The extended specimen:
497      Emerging frontiers in collections-based ornithological research, 1–9. CRC Press,
498      Boca Raton, Florida, USA.
499  Wilf P., Zhang S, Chikkerur S, Little SA, Wing SL, Serre T. 2016. Computer vision
500      cracks the leaf code. *PNAS* 113: 3305–3310.
501  Willis CG, Law E, Williams AC, Franzone BF, Bernardos R, Bruno L, Hopkins C,
502      Schorn C, Weber E, Park DS, Davis CC. 2017. CrowdCurio: An online
503      crowdsourcing platform to facilitate climate change studies using herbarium
504      specimens. *New Phytologist* 215: 479–488.
505  Yost JM, Sweeney PW, Gilbert E, Nelson G, Guralnick R, Gallinat AS, Ellwood ER,
506      Rossington N, Willis CG, Blum SD, Walls RL. 2018. Digitization protocol for
507      scoring reproductive phenology from herbarium specimens of seed plants.
508      *Applications in Plant Sciences* 6: e1022.
509  Zhou N, Siegel ZD, Zarecor S, Lee N, Campbell DA, Andorf CM, Nettleton D,
510      Lawrence-Dill CJ, Ganapathysubramanian B, Kelly JW, Friedberg I. 2018.
511      Crowdsourcing image analysis for plant phenomics to generate ground truth data
512      for machine learning. *PLoS Computational Biology* 14: e1006337.
513

514  **Figures.**
515
516  **Figure 1**. Digitization 1.0 and 2.0. Digitization 1.0 is the creation and online mobilization
517  of digital content derived from physical specimens. This endeavor occurs locally within
518  institutions, most commonly Natural History Museums. Digitization 2.0, in contrast,
519  builds upon the digitized data, workflows, and infrastructure produced by Digitization 1.0
520  to facilitate enhanced digitization, curation, and data linkages to address increasingly
521  complex questions at a massive global scale not previously imagined. These efforts are
522  stimulating a new work force and connecting diverse scholarly domains, propelling a new
523  generation of scientists forward, and removing financial, sociological, institutional, and
524  academic obstacles restricting access to these materials. Some areas of inquiry that will
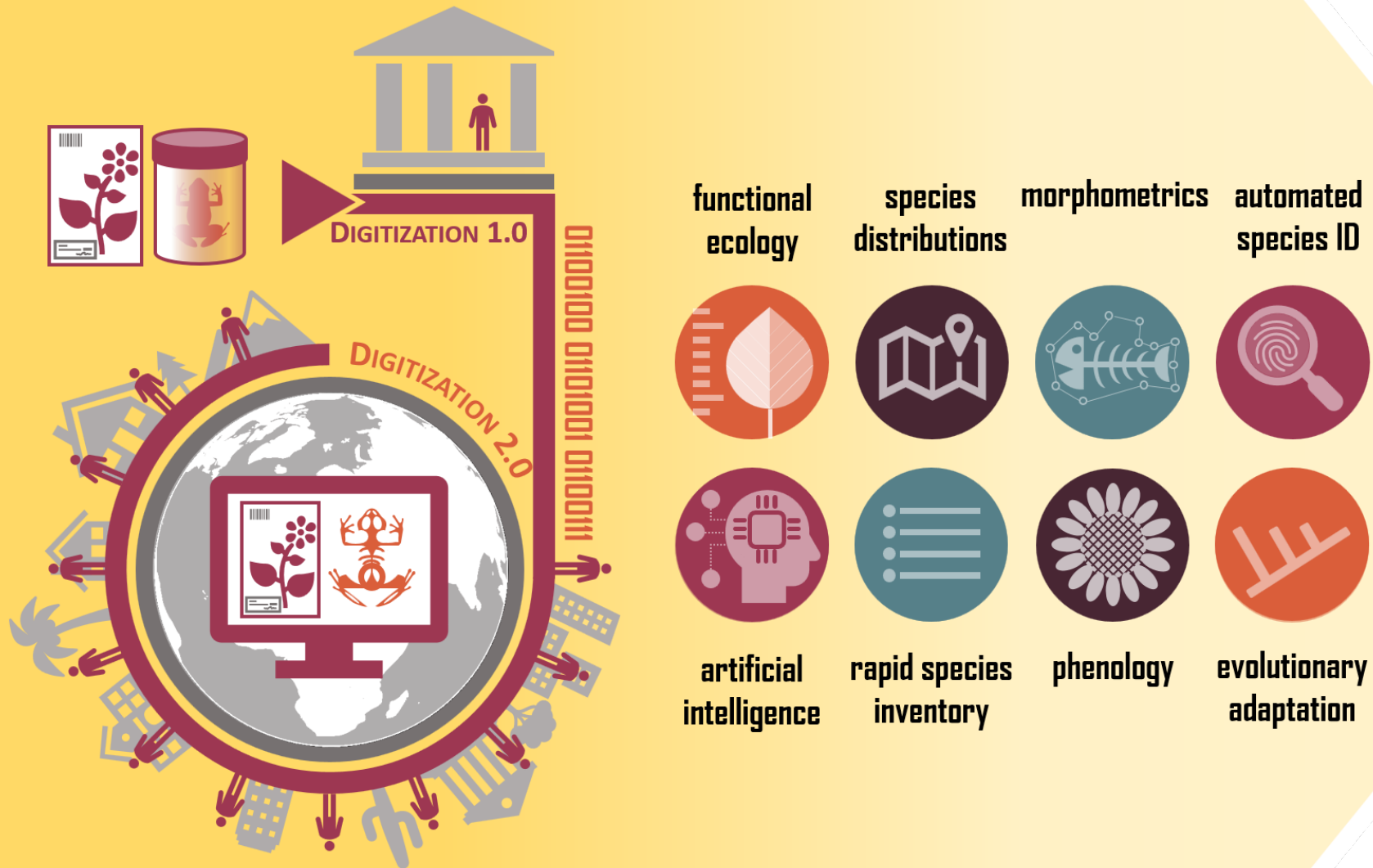525  be greatly stimulated by both Digitization 1.0 and 2.0 are highlighted.
526
527  **Figure 2.** An end-to-end pipeline example to highlight the value and complementarity of
528  Digitization 1.0 and 2.0. The African pig-nosed frog (genus *Hemisus*) shown (A) was
529  collected during recent field research in Angola. In addition to metadata from the
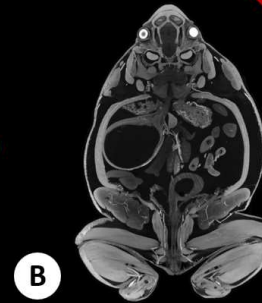530  collection event, a series of x-ray images (tomograms) were created using diffusible

531    iodine-based contrast-enhanced computed tomography (diceCT) directly from the
532    voucher specimen. This product of Digitization 1.0 is shown in the black and white x-ray
533    image (B). From these digital x-ray images, a 3D volume was created from the digital
534    data generated during Digitization 1.0 from which students and scientists can digitally
535    dissect and manipulate regions of interest representing the frog's nervous (C), circulatory
536    (D), and muscular (E) systems (Digitization 2.0).
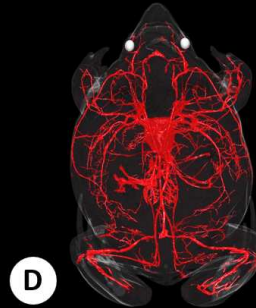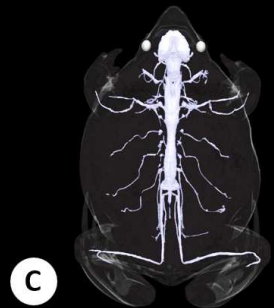537
538    **Figure 3.** Estimating collection sizes and impact on research. (A) Size and geographical
539    distribution of the vascular plant collection at the Harvard University Herbaria (HUH).
540    To statistically estimate the size of this large collection, the total number of specimens in
541    randomly subsampled cubbies were counted. These data were then used to model a
542    probability distribution of the total number of specimens across the entire collection
543    (Comoglio et al. 2013). Three hundred fifty cubbies were sampled and counted,
544    establishing that the HUH has 3,701,695 vascular plants with a 95% confidence interval
545    spanning 3,644,497 to 3,759,803. A similar approach was applied to further assess
546    geographical distribution of the collection as well. (B) Loan use information for the
547    Harvard Museum of Comparative Zoology ichthyology collection. Digitization greatly
548    enhances the tracking of loan use history post 1980, until which point records are
549    confined to physical logbooks. (C) Cumulative number of HUH specimen loans post
550    1980. While the total number of physical specimen loans (red) have remained relatively
551    constant in recent years, the number of digital specimen images loaned has grown
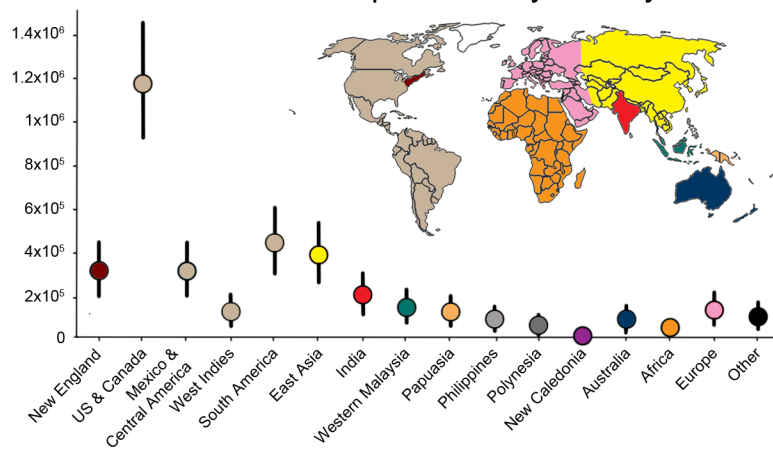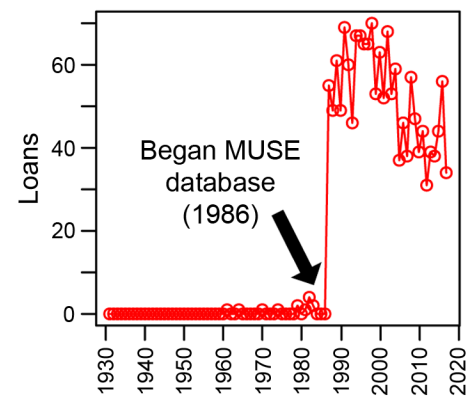552    substantially.

Digitization 1.0

Digitization 2.0

functional ecology

species distributions

morphometrics

automated species ID

artificial intelligence

rapid species inventory

phenology

evolutionary adaptation

DIGITIZATION 1.0

DIGITIZATION 2.0

01100100 01101001
01100111 01101001
01110100 01101001
01111010 01100001
01110100 01101001
01101111 01101110

**A** Vascular Plant Specimens by Locality

**B** MCZ Ichthyology Loans

Began MUSE database (1986)

**C** First Digital Loans