

Digitization and the Future of Natural History Collections

BRANDON P. HEDRICK, J. MASON HEBERLING, EMILY K. MEINEKE, KATHRYN G. TURNER, CHRISTOPHER J. GRASSA, DANIEL S. PARK, JONATHAN KENNEDY, JULIA A. CLARKE, JOSEPH A. COOK, DAVID C. BLACKBURN, SCOTT V. EDWARDS, AND CHARLES C. DAVIS

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 on the data, workflows, and infrastructure produced by Digitization 1.0 to create digital-only workflows that facilitate digitization, curation, and data links, 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 issues of global change.

Keywords: digitization, herbaria, natural history collections, specimens, Anthropocene, baselines

Anthropogenic impacts, including urbanization, globalization, and climate change, are rapidly transforming our world. Despite our best efforts, however, quantifying the biotic impacts of human activity has been challenging, as is evidenced by the difficulty of delimiting the onset of the Anthropocene (Lewis and Maslin 2015). Part of this uncertainty stems from a lack of historical data that track biotic change through time. However, natural history collections (NHCs), with their broad taxonomic, geographic, and temporal scope, offer a key solution to this impasse. In the past 20 years, there has been a dramatic increase in the use of NHCs for assessing a wide variety of scientific questions (Suarez and Tsutsui 2004, Pyke and Ehrlich 2010, Park and Potter 2015, Meineke et al. 2018, 2019). Indeed, they have emerged as one of the best resources for establishing biological baselines to understand the impacts of, for example, the origins of agriculture, the industrial revolution, the development of nuclear armaments, and—more generally—the influence and acceleration of anthropogenic change on biodiversity (Moritz et al. 2008, Johnson et al. 2011, Lister 2011, Funk 2018, Nelson and Ellis 2018).

Most large NHCs provide specimen data to researchers and the public by mobilizing searchable collection databases online. We assert that these mobilized collections are among the most important advances in museum curation in the past century, significantly opening access to NHCs and greatly stimulating large-scale analyses that span novel academic and

societal enterprises. These resources are connecting diverse scholarly domains, propelling a new generation of scientists forward, and removing financial, sociological, institutional, and academic obstacles preventing access to these materials (Drew et al. 2017, Sweeney et al. 2018). In short, digitizing a specimen—translating metadata associated with a physical specimen object into flexible digital data formats—increases the value of the physical specimen exponentially.

In the present article, we present an ambitious, two-pronged vision for digitization, which we term Digitization 1.0 and Digitization 2.0. Digitization 1.0 represents the ongoing push to create digital images and related content directly from physical voucher specimens; Digitization 2.0, in contrast, relates exclusively to data gathering, tasks, or workflows derived from digitized products of Digitization 1.0 rather than from the physical specimens themselves (figure 1). In addition to the vast expansion and online aggregation of these mobilized collections to create a truly global digital NHC, Digitization 2.0 also expands the process of digitization globally and expands the workforce that interfaces with these objects by including researchers far from NHCs and also citizen scientists, thus serving to accelerate the progress of digitization.

Digitization 1.0: The past, present, and future

The digitization of NHCs began with the overarching goal of documenting specimen inventories and facilitating research by transcribing label information into centralized,

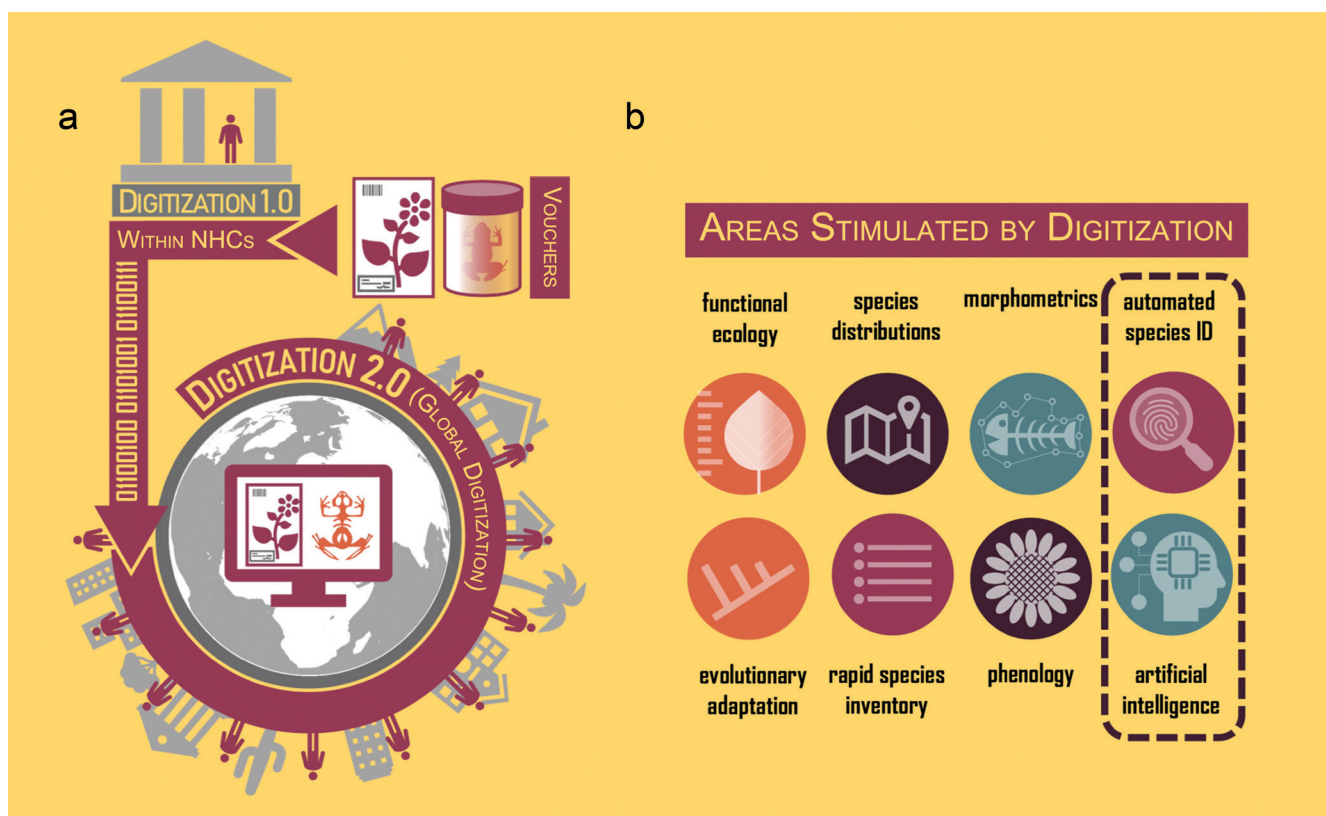


Figure 1. Digitization 1.0 and 2.0. Digitization 1.0 is the creation and online mobilization of digital content derived from physical specimens. This endeavor occurs locally within institutions, most commonly Natural History Museums. Digitization 2.0, in contrast, builds on the digitized data, workflows, and infrastructure produced by Digitization 1.0 to facilitate enhanced digitization, curation, and data links to address increasingly complex questions at a massive global scale not previously imagined. These efforts are stimulating a new work force and connecting diverse scholarly domains, propelling a new generation of scientists forward, and removing financial, sociological, institutional, and academic obstacles restricting access to these materials. (a) A simplified digitization pipeline where voucher specimens held in NHCs are digitized within NHCs themselves (Digitization 1.0) and then those digitized products are distributed globally for additional digitization by researchers, citizen scientists, and through automated digitization (Digitization 2.0). (b) Some areas of inquiry that have been greatly stimulated by both Digitization 1.0 and 2.0 are highlighted. Areas stimulated only through Digitization 2.0 are bounded by a black outline.

searchable databases, as was described recently by Nelson and Ellis (2018). These efforts have given rise to Digitization 1.0, which has been widely embraced and continues to be infused with innovation. Digital representations generated through Digitization 1.0 include specimen images and direct transcriptions of specimen metadata from handwritten or printed collection catalogs or labels, including, for example, details on coloration or measurements. As part of this effort, NHCs have generated millions of digital representations of physical vouchers and have devised numerous technological innovations to facilitate efficient data generation, including conveyor belt and robotic imaging techniques for mass specimen digitization (Tegelberg et al. 2014, Sweeney et al. 2018). More recent next-generation technologies, including photogrammetry, laser scanning, and computed tomography, create far richer digital representations of specimens than can

be visualized by eye or with standard microscopy (figure 2). Given that large portions of most NHCs still remain unavailable in digital format, the innovations and efforts within Digitization 1.0 will continue well into the future, likely for decades. In the subsections below, we outline Digitization 1.0 through the lens of digitization workflows, strategic prioritization, and solutions to impediments.

Digitization workflows and linking data. The practice of digitization is broadly consistent among projects and organismal groups, in so much as each specimen is represented by textual metadata from labels or catalogs and typically two-dimensional digital images but increasingly also three-dimensional representations and audio or video recordings where relevant. There exists great variation in specimen size, storage conditions (e.g., fluid preserved, microscope slides,

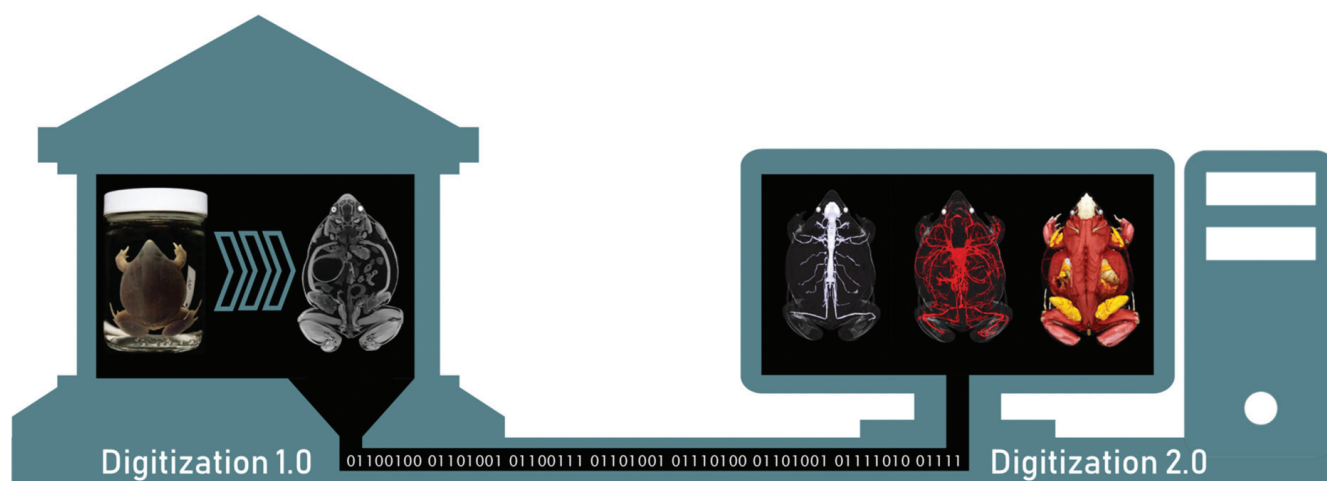


Figure 2. An example pipeline to highlight the value and complementarity of Digitization 1.0 and 2.0. The African pig-nosed frog (genus *Hemisus*) shown was collected during recent field research in Angola. In addition to metadata from the collection event, a series of X-ray images (tomograms) were created using diffusible iodine-based contrast-enhanced computed tomography (diceCT) directly from the voucher specimen within an NHC (Digitization 1.0). This three-dimensional digital data (e.g., CT data) generated during Digitization 1.0 can be digitally dissected and manipulated by researchers or citizen scientists anywhere on the globe to highlight the frog's nervous, circulatory, and muscular systems (Digitization 2.0).

dry storage), dimensionality (two- versus three-dimensional representation), and detail associated with specimens, not to mention widely varying practices in specimen collection and curation across taxonomic domains and institutions. This heterogeneity of collections and institutional policies and priorities therefore creates challenges to efficient mass imaging and harvesting of metadata. However, at minimum, digitization workflows should attempt to integrate all available specimen metadata into digitization efforts and appropriately link these data to their associated physical voucher specimens. Beyond traditional linkages, nontraditional metadata associated with specimens include biotic (e.g., mass) and abiotic data (e.g., climate), media (e.g., video and audio recordings), community- and population-level metadata (e.g., abundance), species observations in the field, and genetic samples (i.e., the extended specimen; Lendemer et al. 2019). Many of these digital data are served in part or in their entirety via online collection database platforms and management software (e.g., Arctos, Specify, Symbiota, EMu) or in data aggregators (e.g., iDigBio, Global Biodiversity Information Facility [GBIF], Botanical Information and Ecology Network). Linking voucher specimens to these new data layers post collection is important and has been facilitated by associating URLs, data accession numbers, DOIs (digital object identifiers), or ARKs (archival research keys) with specimen records in collection databases. In addition, trait data can be incorporated into specimen records using extensions to the Darwin Core Archives (Yost et al. 2018). For the next generation of collections, protocols are under development to expand the digitization workflow to the collecting event itself (Heberling and Issac 2018).

Developing digitization priorities. Given the limited resources available to many NHCs, it is necessary to establish priorities for specimen digitization. Specimens at risk of degradation, such as rare or fragile fossils, and those representing rare or threatened species and habitats are candidates for high priority digitization. Furthermore, efforts should focus on specimens with rich associated metadata from the collection event. A growing number of species are imperiled, and conservation biologists are increasingly reliant on NHCs for baseline data to understand species ranges and climatic tolerances for assessing future changes (Lister 2011). Distributing information for these rare or threatened taxa to conservation biologists is increasingly critical to these species' management and survival (MacDougall et al. 1998, Nualart et al. 2017). Finally, taxa representing either a breadth of evolutionary history or unique adaptations are important for research on phenotypic evolution, community ecology, and biologically inspired design. We suggest that such specimens have high priority for digitization.

When attempting to isolate specimens that are of high priority, especially in large collections, it is important to remember that NHCs are assembled nonrandomly, often driven by convenience or opportunity, individual interests, funding priorities, and other logistical factors (Pyke and Ehrlich 2010, Daru et al. 2018). Therefore, certain areas, times, and taxa are much better represented than others, leading to gaps in our knowledge of global biodiversity. These biases are further compounded in ways that are not yet well characterized by disparate efforts to generate digital derivatives of these primary biodiversity data (e.g., arthropod collections are under digitized, despite their scope;

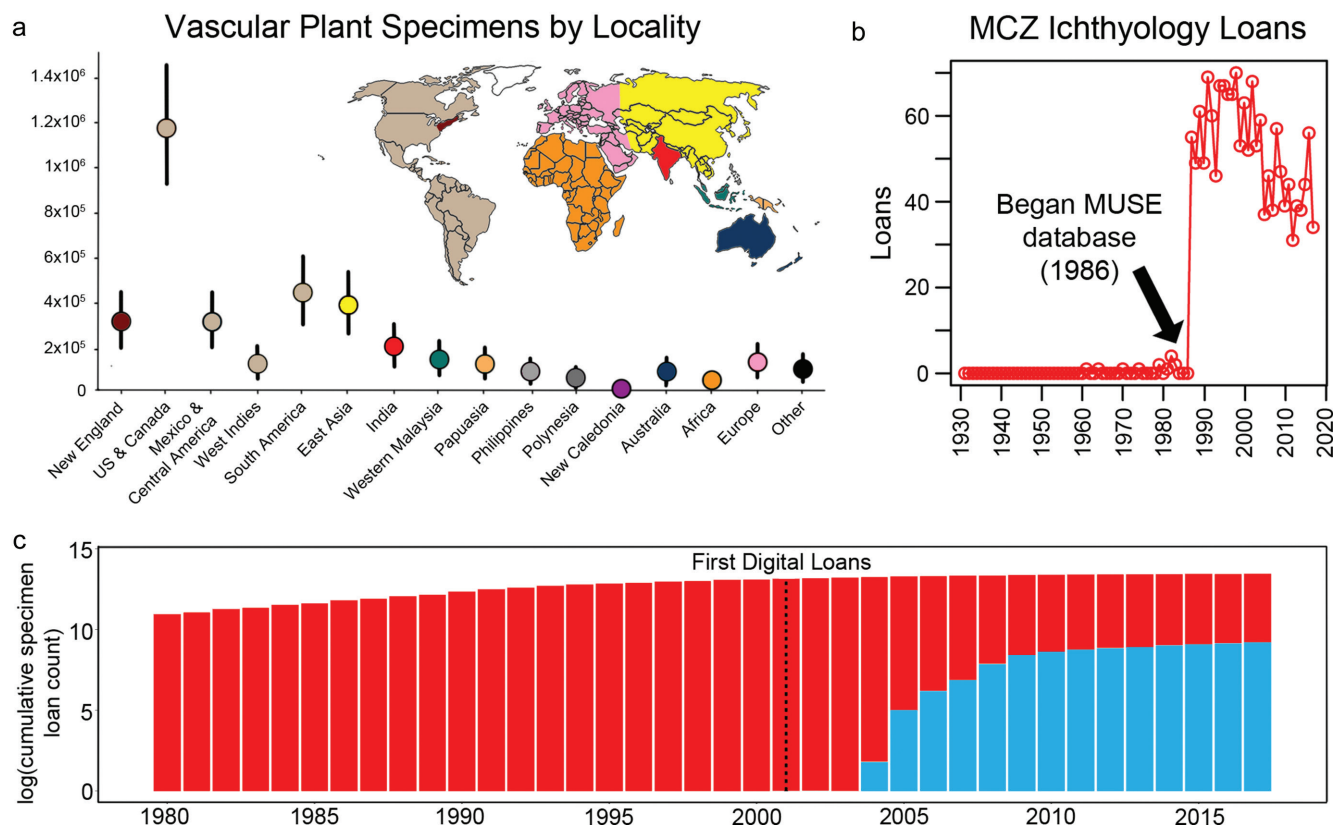


Figure 3. Estimating collection sizes and impact on research. (a) Size and geographical distribution of the vascular plant collection at the Harvard University Herbaria (HUH) showing that the HUH is a representative example of a large NHC with a global distribution of specimens. To statistically estimate the size of this large collection, the total number of specimens in randomly subsampled cubbies were counted. These data were then used to model a probability distribution of the total number of specimens across the entire collection (Comoglio et al. 2013). Three hundred fifty cubbies were sampled and counted, establishing that the HUH has 3,701,695 vascular plants with a 95% confidence interval spanning 3,644,497 to 3,759,803 (see box 1). A similar approach was applied to further assess geographical distribution of the collection as well. (b) Loan use information for the Harvard Museum of Comparative Zoology ichthyology collection. Digitization greatly enhances the tracking of loan use history post 1980, until which point records are confined to physical logbooks. (c) Cumulative number of HUH specimen loans post 1980. Although the total number of physical specimen loans (red) have remained relatively constant in recent years, the number of digital specimen images loaned has grown substantially.

Cobb et al. 2019). To alleviate at least some of the latter bias during the digitization process, comprehensive assessments of spatial, temporal, and taxonomic biases within collections can be used to identify gaps that can be remedied via targeted digitization efforts (Beck et al. 2013, Beck et al. 2014, Meyer et al. 2016, Troudet et al. 2017, Daru et al. 2018).

Owing to the varying effort required by different digitization strategies (e.g., label data, images, three-dimensional reconstructions), data types that serve the largest diversity of use cases should also be prioritized. For instance, key information including taxon name, collection locality, and date can be captured relatively efficiently and can facilitate assessments of species distributions through time. Rapidly expanding areas of research, including phenology (e.g., Primack et al. 2004, Willis et al. 2017), large-scale taxonomic inventories (e.g., Cardoso et al. 2017), and morphometric

investigations (e.g., Hedrick et al. 2015), rely on such label data and data from postdigitization enhancement (Sweeney et al. 2018). For example, in one of the first studies to demonstrate how historic specimens can be used to quantify the biotic effects of climate change, Primack and colleagues (2004) used flowering plant specimens collected between 1885 and 2003 in the greater Boston area to demonstrate that plants were flowering up to 8 days earlier in recent years than in the early years of the twentieth century. The utility of such diverse data (e.g., geographic location, flowering date, anatomical measurements) is important to a wide array of researchers and should be prioritized. In addition, we feel it is best to only apply more complex, holistic digitization methods on a key subset of data-rich specimens as has been recently demonstrated in the openVertebrate (oVert) thematic collection network (Blackburn et al., NSF abstract

Box 1. Estimating the size and scale of a global digitization effort.

Digitization 1.0 has resulted in the mobilization of millions of specimen records and has created the momentum for a massive, global digitization effort. To better establish target goals and evaluate the success of this effort (e.g., estimating the proportion of specimen records that have been digitized and mobilized online), obtaining accurate estimates of the number of specimens housed in NHCs is necessary. Extrapolations from digitized content indicate that roughly 2.5 billion–3 billion specimens are housed in NHCs worldwide (O’Connell et al. 2004, Krishnan et al. 2016). However, more robust assessments of global specimen numbers, including geographic and taxonomic distribution, are required to facilitate thoughtful assessments of collection bias to better target digitization priorities (Meyer et al. 2016). Making robust size estimates are particularly relevant as vended solutions are used to achieve digitization milestones, including the mobilization of entire collections, such as those at the Muséum National D’Histoire Naturelle (France), Naturalis (Netherlands), and the Smithsonian Institution (United States; Rogers 2016, Le Bras et al. 2017). Along these lines, a test case example to illustrate such an effort on a smaller scale comes from the Harvard University Herbaria (HUH), which has been thought to contain 5.5 million specimens. Targeted subsampling of the HUH vascular plant collection facilitated accurate estimates (with confidence intervals) of total specimen collection numbers and their geographic distribution (figure 3a). Once the total number of specimens in NHCs have been accurately quantified, it is necessary to establish the percentage of specimen collection records that have been digitally mobilized.

Because imaging and serving metadata-rich collection information online requires a large financial investment, as well as human labor, its impacts on research should be documented and acknowledged. The most powerful outcomes of digitization would be better characterized by relating these various forms of data usage to one another to explore how digitization increases specimen usage. Along these lines, data relevant to describing the scientific impact of physical specimens (predigitization), such as loans and museum visits, remain largely confined to physical collection logbooks, thus limiting assessment of the impact of Digitization 1.0 (figure 3b). Such efforts would allow us to begin to understand the ways that digitization stimulates increased visitation and use of the actual physical versus digital collection (figure 3c). As a community, we must be better prepared to track and assess these questions.

no.1701714). Increasing the magnitude of the collection of media files (e.g., photogrammetry of bird skins, nuts) for this subset of data via new pipelines and technological advancements will be critical to this effort.

Past impediments and future solutions. Despite the success of Digitization 1.0, this initiative has identified three issues that must be addressed to maximize efficiency of information retention and distribution. First, museums are obligated to manage, store, and steward additional digital data associated with their physical collections. However, the act of digitization entails significant challenges, because it requires sustainably curating both the physical objects and rapidly emerging digital data sets. This issue will necessitate the development of new tools, will require that centralized aggregators assume more responsibility, and will require increased funding in the near future (see the “Digitization 2.0” section below).

Second, there is concern that large aggregators aimed at connecting researchers with NHCs (e.g., GBIF, iDigBio; Edwards 2004) remove NHCs from the attribution chain. NHCs are frequently funded on their research relevance. When researchers view specimen images or harvest metadata from aggregators, NHCs that contribute these data often receive little to no credit (Rouhan et al. 2017). A mechanism for referencing these source collections needs to be embedded in the publication process that requires that NHCs be acknowledged and notified when publications incorporate their data. A viable solution to this problem is to mint a DOI for a digitized specimen and establish a reporting mechanism for collections to be alerted when

their specimens have been cited. Automating this attribution pipeline as part of the digitization workflow better ensures that NHCs receive credit for stewarding both voucher specimens and also digitized data, which is likely to stimulate NHCs to embrace open-access policies for their data.

Third, digitized data are inconsistently and redundantly spread across multiple databases at different scales. NHCs often have their own databases, but some data are additionally deposited in regional databases, taxon-specific databases, and national and international data aggregators. This data dispersion causes information to be input or archived redundantly such that each database has a variant of the postdigitization metadata, leading aggregators to archive either inconsistent or duplicated copies of the same primary data. This problem can be partially circumvented by more communication among data aggregators, as well as between NHCs and aggregators. Algorithms linking specimen numbers between aggregators could ensure that postdigitization enhancement metadata are transferred to all aggregators mentioning particular specimens by unique identifiers such as the specimen-based occurrenceID. This is done internally at iDigBio via the iDigBio Record API, which retains current and previous iterations of a specimen’s data.

Digitization 2.0: Charting a road map for the future

Unlike Digitization 1.0, which directly uses the physical specimen, Digitization 2.0 instead uses the digitized product from Digitization 1.0 for generating additional data and metadata (figure 1). Digitization 2.0 is powered by the online aggregation of these resources and enables digitization to assume new forms and engage vast new workforces. As we

outline below, Digitization 2.0 is already well underway and holds tremendous promise. It includes semi- or fully automated data recording from digitized specimens, which stimulates research and returns value to the physical specimen. In addition, Digitization 2.0 involves an expansion in the workforce engaged in collections science. Finally, Digitization 2.0 leverages NHC resources to create trait databases, either from aggregating and better indexing existing metadata or by allowing researchers or citizen scientists to associate trait annotations with images served from NHC databases.

Innovative tools for automating digitization: Machine learning and neural networks. Given the massive number of specimen images in digital databases with minimal data, an important first step is to better automate data transcription to augment these skeletal records. The enormity of this task is rapidly becoming impossibly large for collections staff to manage without automation. In recent years, machine learning applications utilizing convolutional neural networks have achieved stunning levels of performance in computer vision tasks including image detection and classification (Sudholt and Fink 2016). Neural networks have previously demonstrated promising results for handwriting recognition systems, which could easily be applied to automated label transcription. These forms of innovative technology, which have been applied to medical diagnoses, speech recognition, and driverless cars, are now permeating NHCs (Schuettpelz et al. 2017) and are likely to be enormously useful when trained on existing databases of handwriting samples (Krishnan et al. 2016), as well as those from transcribed labels generated through Digitization 1.0. These models can be further trained using existing semantic field constraints to much more effectively parse specimen metadata into appropriate database fields. Beyond capturing essential minimal data records in an automated manner, neural networks have recently been implemented to accomplish far more sophisticated tasks than text transcription (Wilf et al. 2016, Schuettpelz et al. 2017, Funk 2018). Wilf and colleagues (2016), for example, used computer vision to classify fossil leaf images on the basis of leaf shape and venation with high accuracy. This proved not only to be an efficient protocol for classifying images, but also discovered previously unidentified morphological landmarks potentially useful for species identification and for understanding important evolutionary and ecological innovations not previously documented. The community is now ready for deeper exploration of minimal metadata capture using semi- to total automation.

Furthermore, the declining number of taxonomists in the global workforce severely affects our ability to address key questions concerning biodiversity in the face of global change (Hopkins and Freckleton 2002). Combining taxonomists' expertise (past and present) with student and public training and increased automation will facilitate enhanced specimen curation and will greatly enable biodiversity discovery. Continued robust support for taxonomic research and training is essential. However, given the enormity of the

task at hand and the limited time for this effort, we believe that addressing many taxonomic problems of identification, particularly for well-known groups of organisms, could be greatly facilitated by automation (Dou et al. 2012, Feng et al. 2016, Wäldchen et al. 2018, Kho et al. 2017, Valan et al. 2019). Reasonably successful early efforts are underway to machine learn and automatically identify large subcollections of insects (e.g., butterflies; Schermer and Hogeweg 2018). Although simple taxonomic identification may seem rudimentary, it is the foundation of all biological research, and in many groups remains problematic. For example, it is estimated that more than 50% of tropical plant specimens in NHCs are incorrectly identified (Goodwin et al. 2015). Together with the training of more expert taxonomists and organismal biologists, the widespread use of neural networks to identify specimens and target groups that need attention would enhance collection utility for research, teaching, and management and further motivate the discovery and description of new species.

Expansion of the digitization workforce. Expanding digitization to involve a global workforce is now possible and is stimulated by the increasingly global accessibility of NHCs. This is a major advancement of Digitization 2.0. These new workforces can be developed to supplement existing NHC staff, especially whereby new workforces further digitize specimen data (e.g., transcribing label data) from the millions of specimen images residing in databases that have limited associated metadata. One obvious group to engage in this effort is citizen scientists. NHCs associated with museums typically have departments devoted to public outreach, which can easily be tapped for aid, helping collections staff with the task of digitization while simultaneously providing the public with ownership and agency. Using citizen science in this manner has been fruitful in numerous contexts including the transcriptions of label data, georeferencing, and physical specimen annotations (Hill et al. 2012, Ballard et al. 2017, Ellwood et al. 2015, 2017). For example, CrowdCurio's Thoreau's Field Notes, an online crowdsourcing platform has successfully facilitated climate change studies from thousands of herbarium specimens utilizing thousands of nonexpert crowdsourcers (Willis et al. 2017, Park et al. 2019). Notes from Nature has been very successful in transcribing museum records since its inception, with 1,201,712 classifications as of 15 September 2019, showing both the power of citizen science involvement in digitization and that many aspects of Digitization 2.0 are already underway. Quality control is always a concern in large-scale citizen science projects (Willis et al. 2017, Zhou et al. 2018) and therefore an easy-to-use graphical user interface clearly demonstrating to the public how and what to digitize will be important (e.g., Notes from Nature). This has been accomplished in several research-based projects (Chang and Alfaro 2016, Cooney et al. 2017, Willis et al. 2017, Park et al. 2019). Increasingly, such citizen science efforts are being supplemented by machine-based learning as well (Unger et al.

2016, Wilf et al. 2016, Schuettpeitz et al. 2017, Loriuel et al. 2019). For instance, crowdsourced data can potentially provide reliable and rapid data for training and testing machine learning models, creating a positive feedback loop propelling digitization forward.

Layers of trait annotations. Traits of organisms are fundamental for documenting biodiversity but also for understanding how organisms evolve and respond to changing environments. Building on investments in creating digital NHCs, there is now increasing demand for creating and associating new trait data layers to these collections. For some taxa, these biological data are already captured in the digitized text of a specimen record (e.g., Darwin Core fields: “organismRemarks”). In mammals and birds, it is common to have measurements on the mass and length of both the whole specimen and parts of the specimen (e.g., testes length, wing length). The aggregation of traits from both the initial collecting event and new annotations will stimulate a wealth of questions and generate a better understanding of global biodiversity through the development of standardized trait vocabularies (Kissling et al. 2018). For example, recently developed data-processing tools for the data aggregator VertNet standardized more than 1.5 million measurements for vertebrates using digital data from collections (Guralnick et al. 2016). Users can now search those specimen records by mass and length, as well as download harmonized trait data associated with individual specimens. The latter allows for new explorations of trait variation within and across species, including spatial and temporal patterns in traits associated with specimens that have collecting dates and georeferenced localities (Riemer et al. 2018). By expanding this framework to annotate traits to specimens and utilizing online platforms for even three-dimensional representations of specimens, NHCs can facilitate the capture of not only simple traits, ranging from specimen length to the presence of a flower, but also more complex traits requiring more sophisticated representation (e.g., virtual automated dissection of the vertebrate nervous system).

Conclusions

Digitization facilitates the democratizing of collection-based research and is essential to establishing and evaluating biological baselines to assess the impacts of climate change, land-use changes, species invasions, and the current mass extinction. It allows for the mining of specimen data in much the same way that we explore organismal genomes. The key to further developing Digitization 1.0 and further establishing Digitization 2.0 lies in building on what the research, funding, and policy communities have learned in the several decades since the initiation of this endeavor. Data-rich NHC specimens are useful and provide unique perspectives on the diversity and distribution of a given taxon. However, if a specimen is not searchable, it will likely not be found or studied despite its potential use. We are already witnessing the fruits of the synergy between Digitization 1.0 and 2.0. Many

of the research pursuits that are benefited by digitization have a long history in the biodiversity sciences, including species distribution modeling, assessments of phenological response, and morphometric studies. However, it is critical to appreciate that digitization is facilitating these fields at a magnitude that was previously impossible and is ushering in new domains (e.g., feature detection using machine learning, automated species identification), new questions, and new audiences that are not yet realized (or even imagined). Only with creativity and improved techniques, including automated and semiautomated methods, a better distributed digitization workload making use of new technologies and workforces, and conscientious attention to the attribution chain, will researchers be best able to track ongoing biodiversity change from all existing data. Moreover, even as new technologies and digitization techniques emerge, we will need to always return to physical specimens, in ways that are unimaginable now, to generate novel data to better understand our changing planet. Although we stress the importance of improved methods and practices for digitization, the active collection and continued curation of physical specimens by expert biologists remains the central pillar supporting advancements in evolutionary biology and conservation represented so importantly by NHCs.

Acknowledgments

We thank the members and organizers of the National Science Foundation (NSF) Postdoctoral Research Fellows in Biology 2017 Symposium for discussions contributing to this manuscript. NSF and the Biodiversity Collections Research Coordinating Network supported the meeting (NSF grants no. 1441785 and no. 1746177). Additional funding from the NSF includes awards no. NSF-DBI EF1208835 and no. NSF-DEB 1754584 to CCD. For collections-based data from the MCZ and HUH and for discussion, we thank Mark Omura, Linda Ford, Brendan Haley, Paul Morris, Rachel Hawkins, José Rosado, Stephanie Pierce, and Andrew Williston. We thank Pam Soltis, James Hanken, and Maureen Kearney for providing thoughts and comments on earlier versions of the manuscript. Finally, we thank Scott Collins (editor), Barbara Thiers (reviewer), and an anonymous reviewer for feedback improving the quality of the manuscript.

Hedrick and Davis are co-corresponding authors, and Heberling, Meineke, and Turner all contributed equally to the manuscript's production.

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.
- Beck J, Ballesteros-Mejia L, Nagel P, Kitching IJ. 2013. Online solutions and the ‘Wallacean shortfall’: what does GBIF contribute to our knowledge of species’ ranges? *Diversity and Distributions* 19: 1043–1050.
- Beck J, Böller M, Erhardt A, Schwanghart W. 2014. Spatial bias in the GBIF database and its effect on modeling species’ geographic distributions. *Ecological Informatics* 19: 10–15.

- Cardoso D, et al. 2017. Amazon plant diversity revealed by a taxonomically verified species list. *Proceedings of the National Academy of Sciences* 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.
- Cobb NS, Gall LF, Zaspel JM, Dowdy NJ, McCabe LM, Akito YY. 2019. Assessment of North American arthropod collections: prospects and challenges for addressing biodiversity research. *PeerJ* 7:e8086.
- Comoglio F, Fracchia L, Rinaldi M. 2013. Bayesian inference from count data using discrete uniform priors. *PLOS ONE* 8 (art. 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.
- Daru BH, et al. 2018. Widespread sampling biases in herbaria revealed from large-scale digitization. *New Phytologist* 217: 939–955.
- 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, Stiasny MLJ. 2017. Digitization of museum collections holds the potential to enhance researcher diversity. *Nature Ecology and Evolution* 1: 1789.
- Edwards JL. 2004. Research and societal benefits of the global biodiversity information facility. *Bioscience* 54: 485–486.
- Ellwood ER, et al. 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.
- Feng L, Bhanu B, Heraty J. 2016. A software system for automated identification and retrieval of moth images based on wing attributes. *Pattern Recognition* 51: 225–241.
- 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, Wiczorek J, LaFrance R, Bloom D, Russell L. 2016. The importance of digitized biocollections as a source of trait data and a new VertNet resource. *The Journal of Biological Databases and Curation*:1–13. <https://doi.org/10.1093/database/baw158>.
- 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, et al. 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, et al. 2011. Climate change and biosphere response: Unlocking the collections vault. *Bioscience* 61: 148–153.
- Kho SJ, Manickam S, Malek S, Mosleh M, Dhillon SK. 2017. Automated plant identification using artificial neural network and support vector machine. *Frontiers in Life Science* 10: 98–207.
- Kissling WD, et al. 2018. Towards global data products of essential biodiversity variables on species traits. *Nature Ecology and 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, et al. 2017. The French Muséum National D'histoire Naturelle vascular plant herbarium collection data set. *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 data sets. *Trends in Ecology and Evolution* 26: 153–154.
- Lendemer J, et al. 2019. The extended specimen network: A strategy to enhance US biodiversity collections, promote research and education. *Bioscience*. doi:10.1093/biosci/biz140.
- Lorieux T, et al. 2019. Toward a large-scale and deep phenological stage annotation of herbarium specimens: Case studies from temperate, tropical, and equatorial floras. *Applications in Plant Sciences* 7: e01233–e01233.
- 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.
- Meineke EK, Davis CC, Davies TJ. 2018. The unrealized potential of herbaria for global change biology. *Ecological Monographs* 88: 505–525.
- Meineke EK, Davies TJ, Daru BH, Davis CC. 2019. Biological collections for understanding biodiversity in the Anthropocene. *Philosophical Transactions of the Royal Society B: Biological Sciences* 374: 20170386.
- Meyer C, Weigelt P, Kreft H. 2016. Multidimensional biases, gaps and uncertainties in global plant occurrence information. *Ecology Letters* 19: 992–1006.
- Moritz C, Patton JL, Conroy CJ, Parra JL, White GC, Beissinger SR. 2008. Impact of a century of climate change on small-mammal communities in Yosemite National Park, USA. *Science* 322: 261–264.
- Nelson G, Ellis S. 2018. The history and impact of digitization and digital data mobilization on biodiversity research. *Philosophical Transactions of the Royal Society B*. 374: 20170391.
- Nualart N, Ibáñez N, Soriano I, López-Pujol J. 2017. Assessing the relevance of herbarium collections as tools for conservation biology. *The Botanical Review* 83: 303–325.
- O'Connell AFJ, Gilbert AT, Hatfield JS. 2004. Contribution of natural history collection data to biodiversity assessment in national parks. *Conservation Biology* 18: 1254–1261.
- Park DS, Potter D. 2015. Why close relatives make bad neighbours: phylogenetic conservatism in niche preferences and dispersal disproves Darwin's naturalization hypothesis in the thistle tribe. *Molecular Ecology* 24: 3181–3193.
- Park DS, Breckheimer I, Williams AC, Law E, Ellison AM, Davis CC. 2019. Herbarium specimens reveal substantial and unexpected variation in phenological sensitivity across the eastern United States. *Philosophical Transactions of the Royal Society B: Biological Sciences* 374: 20170394.
- Primack D, Imbres C, Primack RB, Miller-Rushing AJ, del Tredici P. 2004. Herbarium specimens demonstrate earlier flowering times in response to warming in Boston. *American Journal of Botany* 91: 1260–1264.
- Pyke GH, Ehrlich PR. 2010. Biological collections and ecological/environmental research: A review, some observations and a look to the future. *Biological Reviews* 85: 247–266.
- Riemer K, Guralnick RP, White EP. 2018. No general relationship between mass and temperature in endothermic species. *eLife* 7: e27166.
- Rogers N. 2016. Museum drawers go digital. *Science* 352: 762–765.
- Rouhan G, Dorr LJ, Gautier L, Clerc P, Muller S, Gaudel M. 2017. The time has come for natural history collections to claim co-authorship of research articles. *Taxon* 66: 101–1016.
- Schermer M, Hogeweg L. 2018. Supporting citizen scientists with automatic species identification using deep learning image recognition models. *Biodiversity Information Science and Standards* 2: e25268.
- Schuetzpelz E, Frandsen PB, Dikow RB, Brown A, Orli S, Peters M, Metallo A, Funk VA, Dorr LJ. 2017. Applications of deep convolutional neural networks to digitized natural history collections. *Biodiversity Data Journal* 5: e21139.
- Suarez AV, Tsutsui ND. 2004. The value of museum collections for research collections. *Bioscience* 54: 66–74.
- Sudholt S, Fink GA. 2016. PHOCNet: A deep convolutional neural network for word spotting in handwritten documents. Pages 277–282 in 2016 15th International Conference on Frontiers in Handwriting Recognition: ICFHR 2016. IEEE.
- Sweeney PW, Starly B, Morris PJ, Xu Y, Jones A, Radhakrishnan S, Grassa CJ, Davis CC. 2018. Large-scale digitization of herbarium specimens:

- Development and usage of an automated, high-throughput conveyor system. *Taxon* 67: 165–178.
- Tegelberg R, Mononen T, Saarenmaa H. 2014. High-performance digitization of natural history collections: Automated imaging lines for herbarium and insect specimens. *Taxon* 63: 1307–1313.
- Troudet J, Grandcolas P, Blin A, Vignes-Lebbe R, Legendre F. 2017. Taxonomic bias in biodiversity data and societal preferences. *Scientific Reports* 7: 9132.
- Unger J, Merhof D, Renner S. 2016. Computer vision applied to herbarium specimens of German trees: Testing the future utility of the millions of herbarium specimen images for automated identification. *BMC Evolutionary Biology* 16: 248.
- Valan M, Makonyi K, Maki A, Vondráček D, Ronquist F. 2019. Automated taxonomic identification of insects with expert-level accuracy using effective feature transfer from convolutional networks. *Systematic Biology* 68: 876–895.
- Wäldchen J, Rzanny M, Seeland M, Mäder P. 2018. Automated plant species identification: Trends and future directions. *PLOS Computational Biology* 14 (art. e1005993).
- Wilf P, Zhang S, Chikkerur S, Little SA, Wing SL, Serre T. 2016. Computer vision cracks the leaf code. *Proceedings of the National Academy of Sciences* 113: 3305–3310.
- Willis CG, et al. 2017. CrowdCurio: An online crowdsourcing platform to facilitate climate change studies using herbarium specimens. *New Phytologist* 215: 479–488.
- Yost JM, et al. 2018. Digitization protocol for scoring reproductive phenology from herbarium specimens of seed plants. *Applications in Plant Sciences* 6: e1022.
- Zhou N, et al. 2018. Crowdsourcing image analysis for plant phenomics to generate ground truth data for machine learning. *PLOS Computational Biology* 14 (art. e1006337).

Brandon P. Hedrick (bphedrick1@gmail.com) is affiliated with the Department of Cell Biology and Anatomy at Louisiana State University Health Sciences Center, in New Orleans, Louisiana. Brandon P. Hedrick, Emily K. Meineke, and Scott V. Edwards are affiliated with the Department of Organismal and Evolutionary Biology at Harvard University, in Cambridge, Massachusetts. J. Mason Heberling is affiliated with the Section of Botany, at the Carnegie Museum of Natural History, in Pittsburgh, Pennsylvania. Emily K. Meineke, Christopher J. Grassa, Daniel S. Park, Jonathan Kennedy, and Charles C. Davis (cdavis@oeb.harvard.edu) are affiliated with Harvard University Herbaria, at Harvard University, in Cambridge, Massachusetts. Kathryn G. Turner is affiliated with the Department of Biological Sciences, Idaho State University, in Pocatello. Julia A. Clarke is affiliated with the Jackson School of Geosciences, at the University of Texas at Austin, in Austin, Texas. Joseph A. Cook is affiliated with the Department of Biology at the University of New Mexico, in Albuquerque. David C. Blackburn is affiliated with the Florida Museum of Natural History, at the University of Florida, in Gainesville.