

Conference Abstract

Unearthing the Past for a Sustainable Future: Extracting and transforming data in the Biodiversity Heritage Library for climate action

JJ Dearborn[‡], Mike Lichtenberg[‡], Joel Richard[‡], Joseph deVeer[§], Michael Trizna^{||}, Katie Mika[¶]

[‡] Smithsonian Libraries and Archives, Biodiversity Heritage Library, Washington, D.C., United States of America

[§] Harvard University, Museum of Comparative Zoology, Ernst Mayr Library, Cambridge, MA, United States of America

^{||} Smithsonian Institution, Office of the Chief Information Officer, Data Science Lab, Washington, United States of America

[¶] Harvard Library & Institute for Quantitative Social Science, Cambridge, Massachusetts, United States of America

Corresponding author: JJ Dearborn (dearbornjj@si.edu)

Received: 08 Sep 2023 | Published: 11 Sep 2023

Citation: Dearborn J, Lichtenberg M, Richard J, deVeer J, Trizna M, Mika K (2023) Unearthing the Past for a Sustainable Future: Extracting and transforming data in the Biodiversity Heritage Library for climate action.

Biodiversity Information Science and Standards 7: e112436. <https://doi.org/10.3897/biss.7.112436>

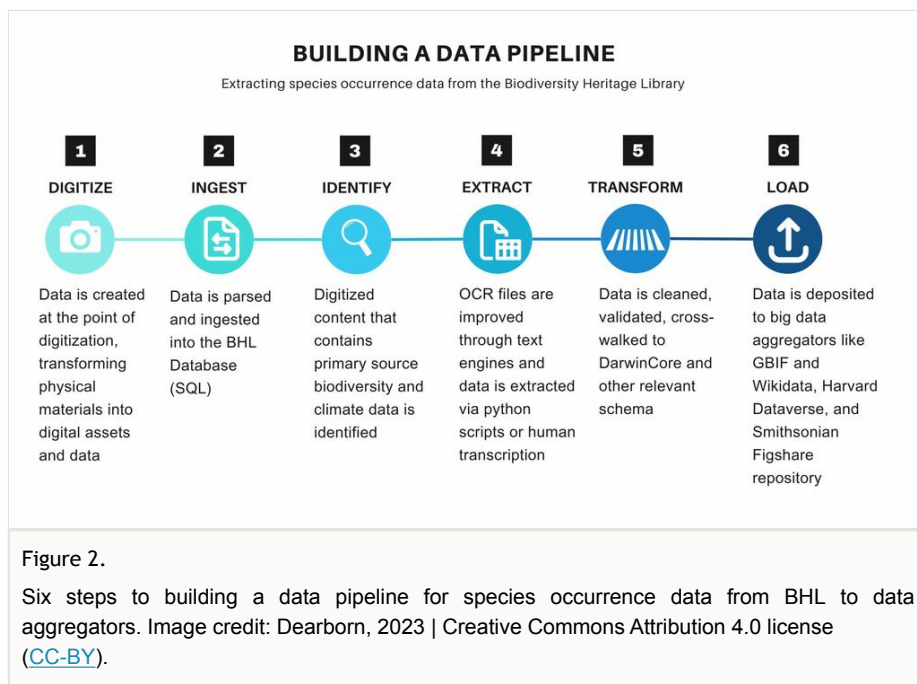
Abstract

As the urgency to address the climate crisis intensifies, the availability of accurate and comprehensive biodiversity data has become crucial for informing climate change studies, tracking key environmental indicators, and building global biodiversity monitoring platforms. The [Biodiversity Heritage Library \(BHL\)](#) plays a vital role in the core biodiversity infrastructure, housing over 60 million pages of digitized literature about life on Earth. Recognizing the value of over 500 years of data in BHL, a global network of BHL staff is working to establish a scalable data pipeline to provide actionable occurrence data from BHL's vast and diverse collections. However, transforming textual content into [FAIR](#) (findable, accessible, interoperable, reusable) data poses challenges due to missing descriptive metadata and error-ridden unstructured outputs from commercial text engines. (Fig. 1)

Despite the wealth of knowledge in BHL now available to global audiences, the underutilization of biodiversity and climate data contained in BHL's textual corpus hinders scientific research, hampers informed decision-making for conservation efforts, and limits our understanding of biodiversity patterns crucial for addressing the climate crisis. By

Textextract) to improve scientific name-finding in BHL's handwritten archival materials using algorithms developed by [Global Names Architecture](#);

- Extracting data from collecting events using HTR coordinate outputs with Python library Pandas DataFrame to create structured data;
- Classifying BHL page-level images with OpenAI's CLIP, a neural network model to accurately identify the handwritten sub-corpus of primary source materials in BHL;
- Running an A/B test to evaluate the efficiency and accuracy of human-keyed transcription data extraction to provide high-quality, human-vetted datasets that can be deposited with data aggregators.



The ongoing development of a scalable data pipeline of BHL's relevant biodiversity and climate-related datasets requires sustained support and partnership with the biodiversity community. Initial results demonstrate that liberating data from archival and handwritten field notes is arduous but feasible. Extending these methodologies to the broader scientific literature presents new research opportunities. Extracting and normalizing data from unstructured textual sources can significantly advance biodiversity research and inform environmental policy. The Biodiversity Heritage Library staff are committed to building multiple scalable data pipelines with the ultimate goal of erecting a global biodiversity knowledge graph, rich in interconnected data and semantic meaning, enabling informed decisions for the preservation and sustainable management of Earth's biodiversity.

Keywords

Global Names Architecture (GNA), global biodiversity infrastructure, scalable data pipelines, historic literature, species occurrence data, climate change, handwritten text recognition, image classification, crowdsourced transcription, structured data, global biodiversity community, data extraction, machine learning, artificial intelligence

Presenting author

JJ Dearborn

Presented at

TDWG 2023

Acknowledgements

Many thanks to additional collaborators from BHL's partner network for making this research possible:

BHL Transcription Upload Tool Working Group (BHL-TUTWG). Members listed in alphabetical order: Riccardo Ferrante (Smithsonian Libraries and Archives), Kelly Hall (Auckland War Memorial Museum, Tamaki Paenga Hira or Auckland Museum), David Iggulden (Kew Gardens Library and Archives), Susan Lynch (BHL), Diane Rielinger (Harvard Botany Libraries), Gretchen Rings (Field Museum, Chicago), Rebekah Kim (California Academy of Sciences), Judy Warnement (BHL).

Global Names Architecture (GNA): Dr. Dmitry Mozzherin, Scientific Informatics Leader at Marine Biological Laboratory and Dr. Geoff Ower, Research Programmer at Illinois Natural History Survey.

Hosting institution

Major support and hosting is provided by Smithsonian Libraries and Archives.

Conflicts of interest

The authors have declared that no competing interests exist.