# BHL, THE BIODIVERSITY HERITAGE LIBRARY:
## Exposing the Taxonomic Literature

Constance Rinaldo, Ernst Mayr Library, Museum of Comparative Zoology, Harvard University, Cambridge MA USA

http://wwwbiodiversitylibrary.org

## WHAT IS THE BHL?

• Large scale digitization to provide open access to core published literature of biodiversity for scientists
• Key component of the *Encyclopedia of Life* http://www.eol.org (EOL) as conceived by E. O. Wilson
• Collaboration of ten major natural history, botanical garden & research libraries: American Museum of Natural History, Field Museum of Natural History, Harvard University (Botany Libraries & Ernst Mayr Library), Marine Biological Laboratory/Woods Hole Oceanographic Institute (MBL/WHOI), Missouri Botanical Garden (MOBOT), Natural History Museum, London, New York Botanical Garden & Royal Botanic Garden, Kew
• Collaboration with global taxonomic community: Global Biodiversity Information Facility (GBIF), International Commission on Zoological Nomenclature (ICZN), European Distributed Institute of Taxonomy, Atlas of Living Australia, Chinese Academy of Sciences, Museum fur Naturkunde der Humboldt-Universitat, BIOONE & more



## WHY DO THIS NOW?

• Biodiversity is HOT; biodiversity studies need taxonomic data
• Taxonomic data are reported in general & specialized literature that may only be found in a few libraries & museums
• Current taxonomic research often relies on multiple texts & specimens more than 100 years old that are dispersed among libraries & museums around the world
• Digital technology offers an access solution to this "taxonomic impediment" that required taxonomists to travel the world to examine every specimen & paper related to an organism
• Taxonomic literature has extreme longevity thus the public domain literature is important
• Literature repatriation: most taxonomic literature is in the developed world while most biodiversity is not (Figure 1)

## WHY A BHL PORTAL?

• Prototype developed at MOBOT as Botanicus.org & tested with scientists
• BHL Portal serves images & text files ingested from Internet Archive
• BHL Portal ingests MARCXML metadata & low resolution JPEG files; High resolution files are retrieved on the fly from IA
• Globally Unique Identifiers (GUIDs) allow links to other services such as EOL
• Taxonomic Intelligence developed at MBL/WHOI allows species name searching by users (Figure 2)
  -TI uses sophisticated algorithm to locate name strings in the Optical Character Recognition (OCR) files that match the 9.4 million names in NameBank
  -Iterative processing of texts increases the number of names in NameBank & the accuracy of recognition
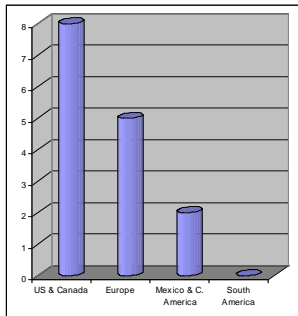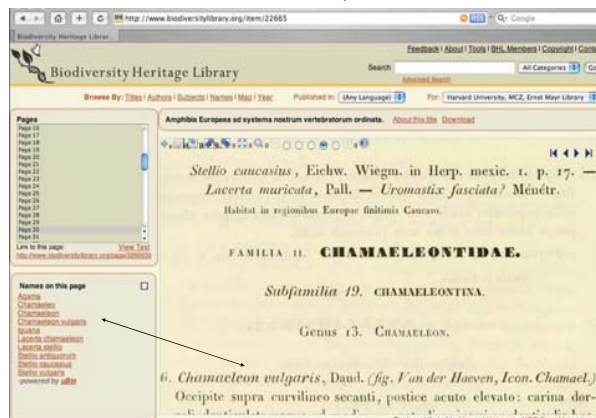  -More tools are under development



Figure 1: Distribution of copies of the *Biologia Centrali-Americana*; the copies in Central America are located in one of the twenty branches of the Smithsonian Libraries. Courtesy, Martin Kalfatovic.



Figure 2: Taxonomic intelligence in action

## WHAT ABOUT COPYRIGHT?

• Public domain literature digitized first
• Opt-in copyright model: BHL actively works with professional societies & other small publishers to integrate publications into the BHL.
• Agreements to digitize 46 titles have been signed with the BHL providing digitization at no cost to society & museum publishers with material served from BHL portal & files available to publishers
• Discussions with commercial publishers for alternative agreements

## HOW?

• BHL is not a legal entity: the ten member institutions signed separate Memoranda of Agreements with the BHL
• Directors of the member libraries meet annually; an elected executive council has weekly conference calls with the BHL Program Director & Technical Director BHL member institution staff have regular conference calls to ensure that all institutions are
• Each institution has a separate contract with Internet Archive, the digitization partner
• IA has small scanning centers in London, DC & Illinois & large centers at the Boston Public Library (thanks to the Boston Library Consortium) & the New York Public Library
• Service is provided for $.10 per page with extra charges for foldouts
• MOBOT, NYBG, Harvard & the Smithsonian have "boutique" scanning facilities to digitize oversized & unusual items
• IA provides image files & text derived from OCR
• OCLC Collection Analysis tool generated a broad look at institutional collection strengths & provided an estimate of the number of public domain materials available for immediate digitization
• Duplication is minimized using tools developed by member libraries such as a serials bidding tool, monograph de-duping tool & others
• Workflow within the libraries includes generating picklists, identifying acceptable items within the picklist, barcoding, generating packing lists, checking out books, packing books, checking in & reshelving returned books & reviewing rejected items

## WHERE DO WE GO FROM HERE?

• Article-level analysis of serials using automated tools
• Further develop global partnerships & incorporate multiple languages
• Linkages to molecular, morphological & other data types
• Improved OCR for non-Roman & non-standard scripts
• Enhance connections with EOL & others
• Expand content access & tools to new audiences
• Strengthen underlying architecture
• Further develop partnerships with commercial & society publishers
• Ingestion of collections that are open access & available

## BIBLIOGRAPHY

Godfray, H.C.J., B.R. Clark, I. J. Kitching, S.J. Mayo & M.J. Scoble. 2007. "The Web and the Structure of Taxonomy," *Systematic Biology*, 56(6): 943-955.
Leary, P.R. , D. P. Remson, C.N. Norton, D.J. Patterson & I.N. Sarkar. 2008. "uBioRSS: Tracking Taxonomic Literature Using RSS," *Bioinformatics* 23(11): 1434-1436.
Minelli, A. 2003 "The Status of Taxonomic Literature," *Trends in Ecology and Evolution* 18(2):75-78.
Sarkar, I.N., R. Schenk & C.N. Norton. 2008. "Exploring Historical Trends Using Taxonomic Name Data," *BMC Evolutionary Biology* 8:144.