

Retriever: Data Retrieval Tool

Henry Kironde¹, Benjamin D. Morris¹, Akash Goel³, Andrew Zhang⁴, Akshay Narasimha⁵, Shivam Negi⁶, David J. Harris⁴, Deborah Gertrude Digges¹, Kapil Kumar⁷, Amritanshu Jain⁵, Kunal Pal⁸, Kevinkumar Amipara⁹, Prabh Simran Singh Baweja¹, and Ethan P. White^{1, 2}

¹ Department of Wildlife Ecology and Conservation, University of Florida ² Informatics Institute, University of Florida ³ Delhi Technological University, Delhi ⁴ The University of Florida ⁵ Birla Institute of Technology and Science, Pilani ⁶ Manipal Institute of Technology, Manipal ⁷ National Institute of Technology, Delhi ⁸ RWTH Aachen University, Aachen, Germany ⁹ Sardar Vallabhbhai National Institute of Technology, Surat

DOI: [10.21105/joss.00451](https://doi.org/10.21105/joss.00451)

Software

- [Review](#) ↗
- [Repository](#) ↗
- [Archive](#) ↗

Licence

Authors of JOSS papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC-BY](#)).

Summary

The Data Retriever automates the first steps in the data analysis workflow by downloading, cleaning, and standardizing tabular datasets, and importing them into relational databases, flat files, or programming languages (Morris and White 2013). The automation of this process reduces the time for a user to get most large datasets up and running by hours to days. The retriever uses a plugin infrastructure for both datasets and storage backends. New datasets that are relatively well structured can be added adding a JSON file following the Frictionless Data tabular data metadata package standard (frictionlessdata 2017). More complex datasets can be added using a Python script to handle complex data cleaning and merging tasks and then defining the metadata associated with the cleaned tables. New storage backends can be added by overloading a general class with details for storing the data in new file formats or database management systems. The retriever has both a Python API and a command line interface. An R package that wraps the command line interface and a Julia package that wraps the Python API are also available.

The 2.0 and 2.1 releases add extensive new functionality. This includes the Python API, the use of the Frictionless Data metadata standard, Python 3 support, JSON and XML backends, and autocompletion for the command line interface.

References

- frictionlessdata. 2017. “Specs: Specifications for Frictionless Data.” <https://github.com/frictionlessdata/specs>.
- Morris, Benjamin D., and Ethan P. White. 2013. “The Ecodata Retriever: Improving Access to Existing Ecological Data.” *PLoS ONE*, June. Public Library of Science (PLOS). doi:10.1371/journal.pone.0065848.