



Collaborative Research: Elements: Advancing Data Science and Analytics for Water (DSAW)

PI: Jeffery S. Horsburgh (jeff.horsburgh@usu.edu)^a, Co-PIs: Brian Crookston^a, Alfonso Torres-Rua^a, Tianfang Xu^b, Anthony Castronova^c

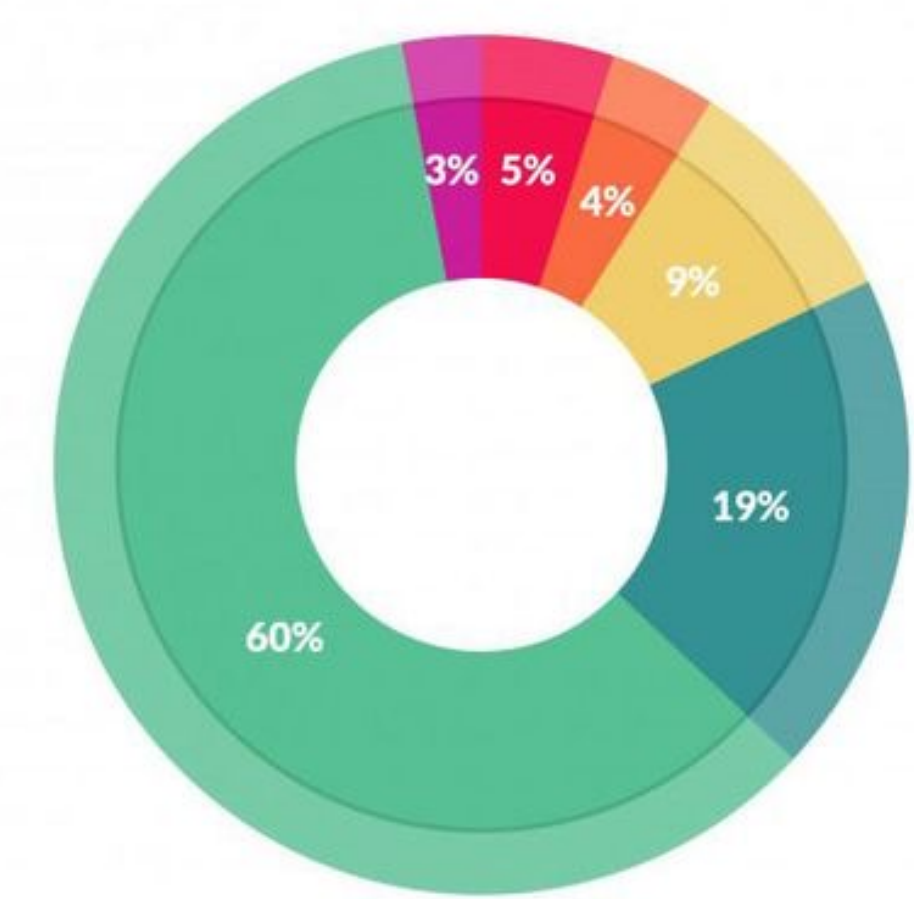
Institutions: ^aUtah Water Research Laboratory, Utah State University; ^bArizona State University; ^cCUAHSI

Award #: 1931297

1. Motivation

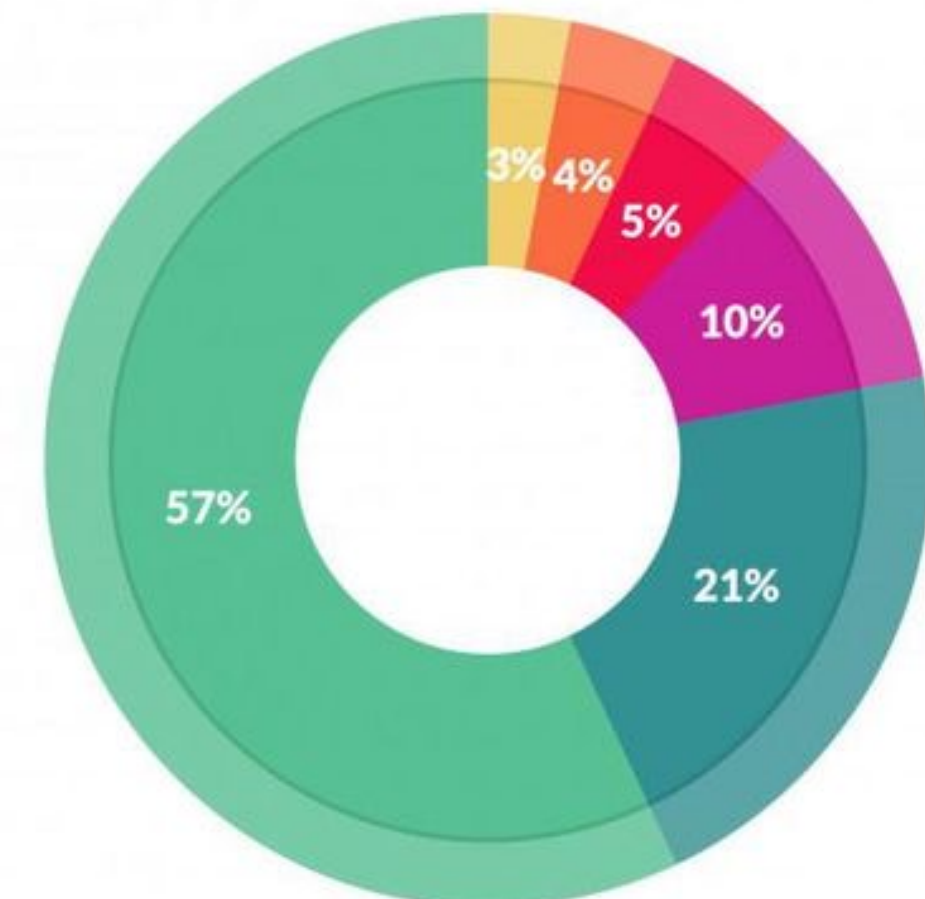
- Challenges in the water domain are multi-disciplinary, requiring data synthesis
- Data manipulation, visualization, and analysis tasks are difficult because datasets are becoming larger, more numerous, and more complex
- Standard data formats for common data types are not always agreed upon or mapped to efficient structures for visualization and/or analysis within an analytical environment.
- Researchers and practitioners lack training in data intensive scientific methods that would enable them to use new and existing data science tools to efficiently tackle large and/or complex datasets
- Overcoming barriers with accessing, organizing, and preparing datasets for data science intensive analyses will be an enabler for transforming scientific inquiries in the water domain.

What do data scientists spend their time doing?



What data scientists spend the most time doing

- Building training sets: 3%
- Cleaning and organizing data: 60%
- Collecting data sets: 19%
- Mining data for patterns: 9%
- Refining algorithms: 4%
- Other: 5%



What's the least enjoyable part of data science?

- Building training sets: 10%
- Cleaning and organizing data: 57%
- Collecting data sets: 21%
- Mining data for patterns: 3%
- Refining algorithms: 4%
- Other: 5%

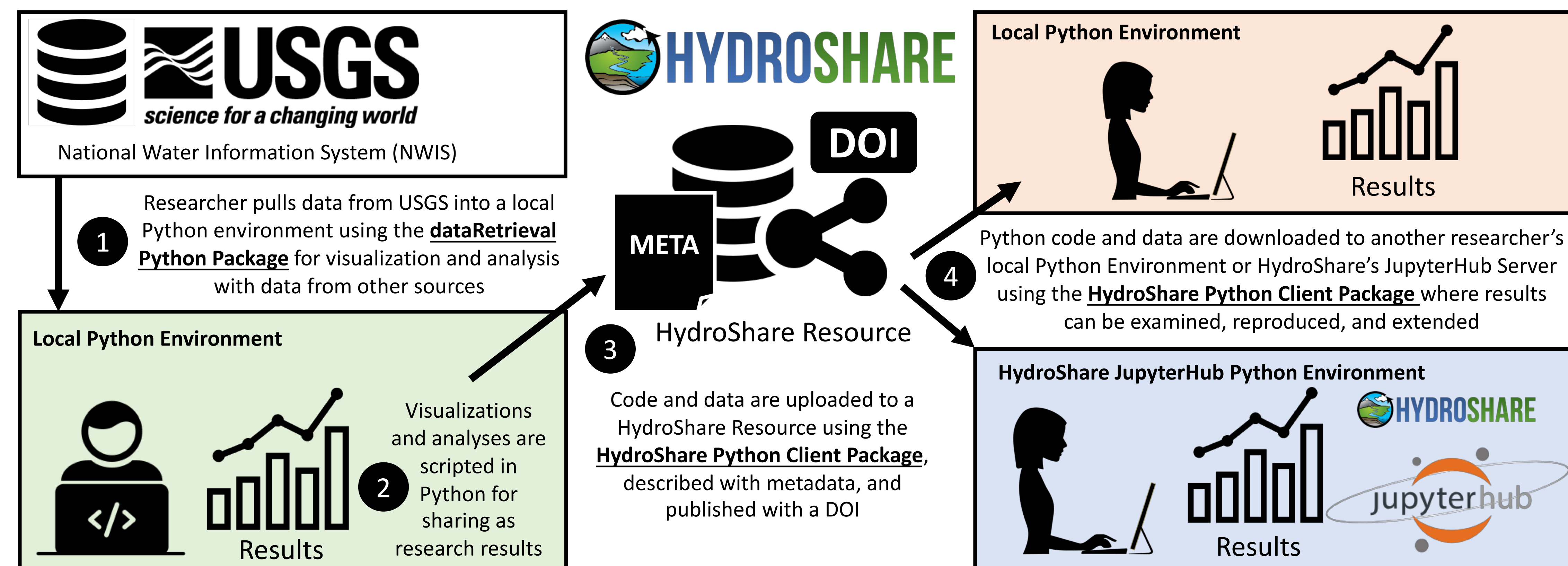
It remains hard to get data and get it into the format you want/need

Cleaning Big Data: Most Time-Consuming, Least Enjoyable Data Science Task, Survey Says

<https://www.forbes.com/sites/gilpress/2016/03/23/data-preparation-most-time-consuming-least-enjoyable-data-science-task-survey-says/?sh=23ad2a316f63>

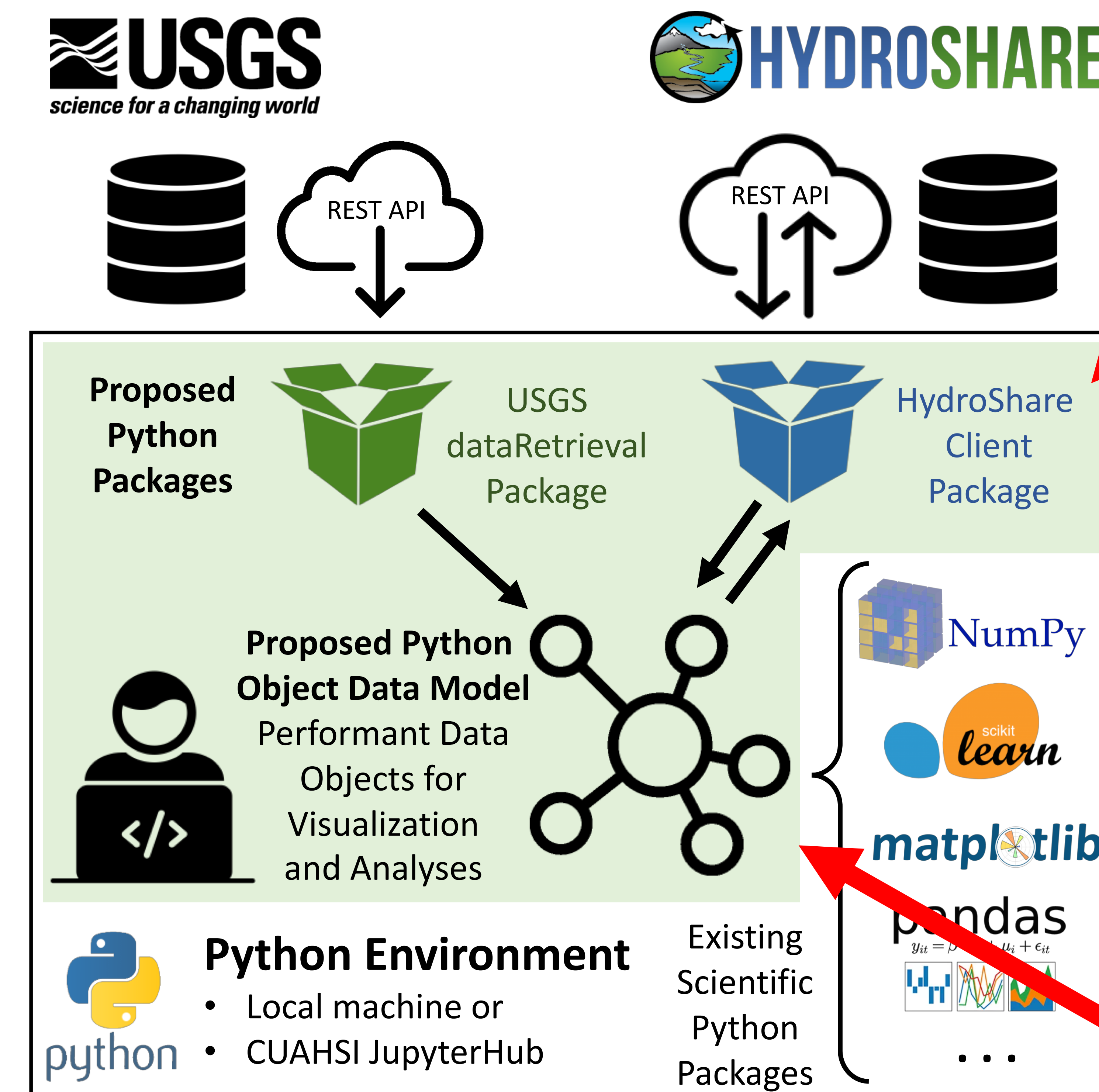
2. DSAW is focused on:

- Enabling water-data scientists to more easily share and collaborate around data and analyses
- Providing data management, visualization, and analysis tools that advance scientists' data science capabilities
- Promoting more consistent data workflows, data reuse, and reproducibility of scientific results.



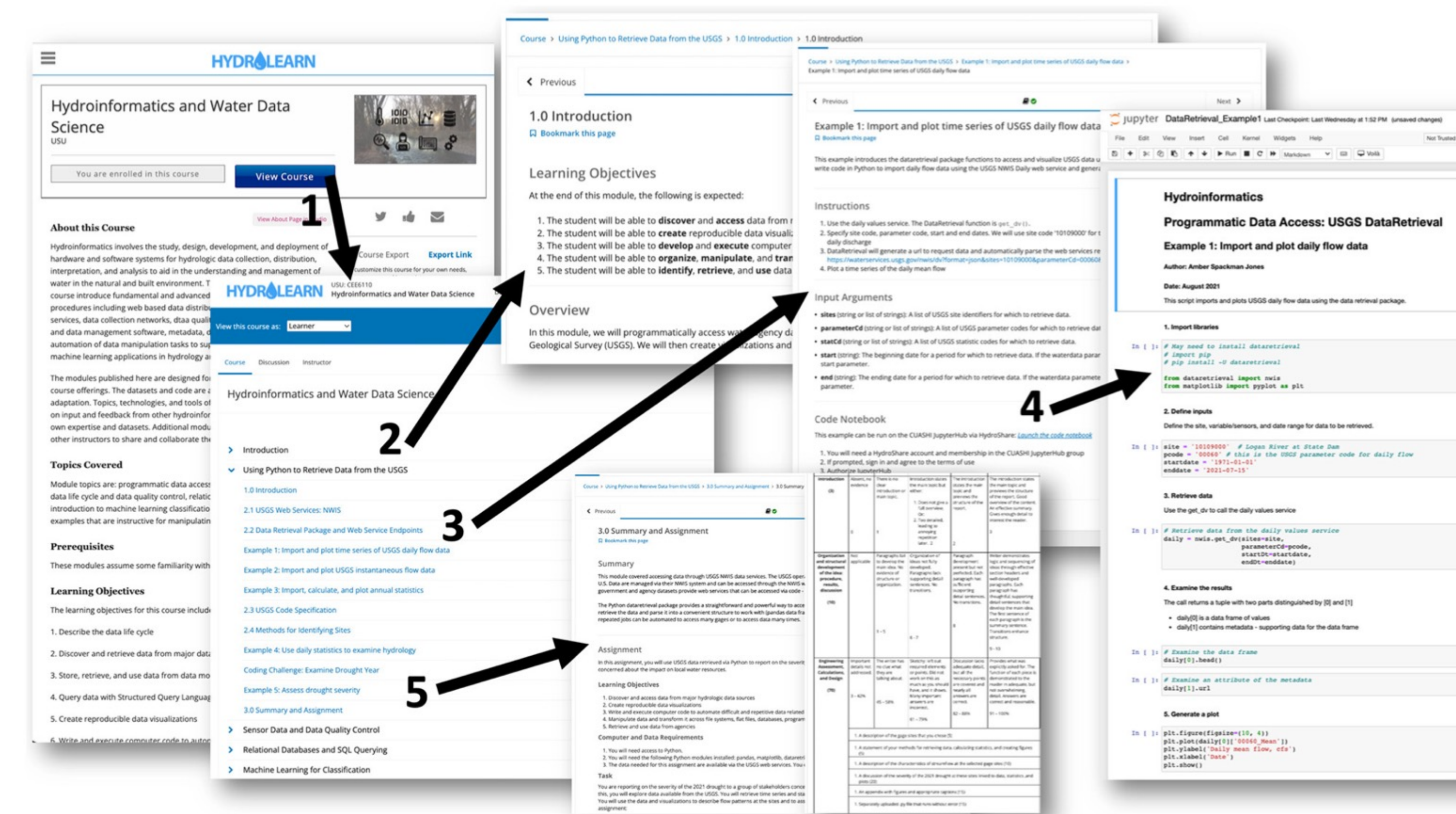
Example reproducible water-data science workflow using HydroShare and new DSAW Python packages.

3. Objectives and Outcomes



DSAW software architecture. New software elements are highlighted in green.

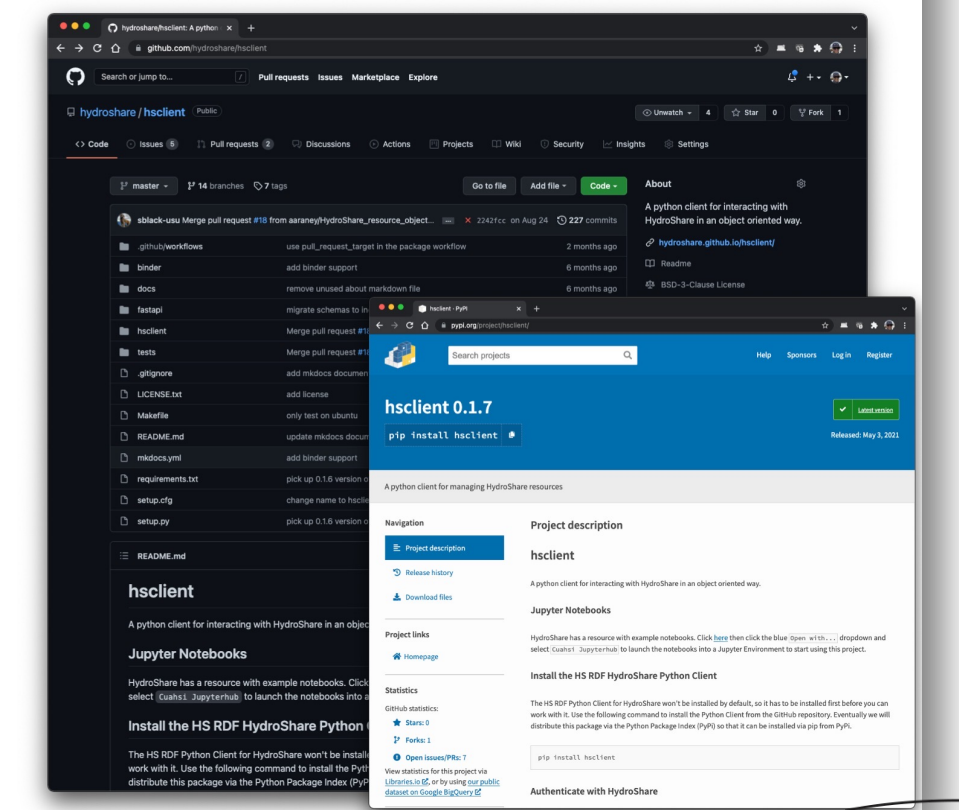
Water-data science applications and educational modules that demonstrate advanced analytical workflows using the tools we develop and help deliver the software to the water science community for use by educators and students



New Python packages that integrate performant data structures with data science capabilities in Python and advanced data access, collaboration, and data archival capabilities of HydroShare

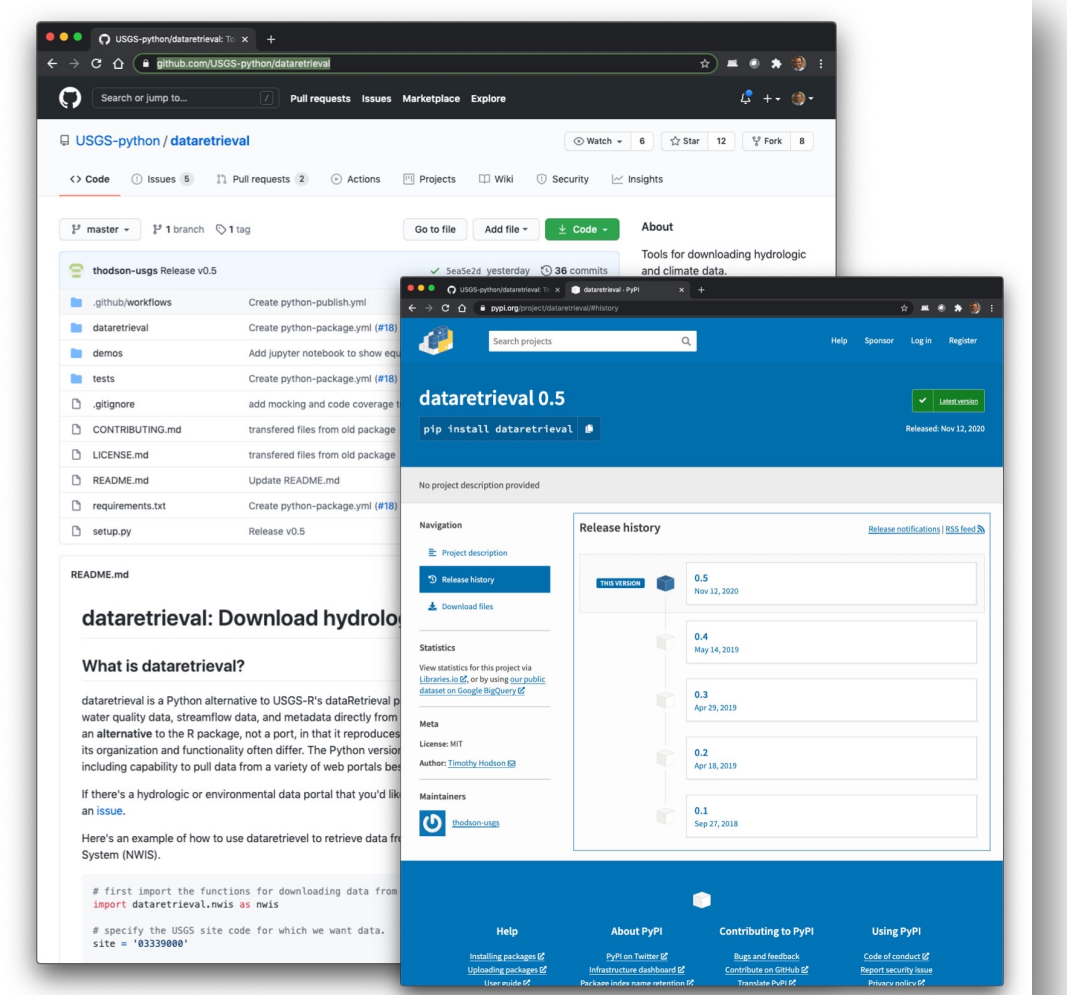
HydroShare Python Client 'hsclient' package

- A set of Python functions for interacting with HydroShare
 - Resource creation/editing
 - Interact with resources in an interactive, object-oriented way
 - Integrate HydroShare resources into data science workflows
 - Reduce the time required to get data for analysis and then save results
- Example Jupyter Notebooks:
<https://www.hydroshare.org/resource/7561a12fd824ebb88edbee05af19b910/>
- GitHub Repository:
<https://github.com/hydroshare/hsclient>



USGS dataretrieval Python package

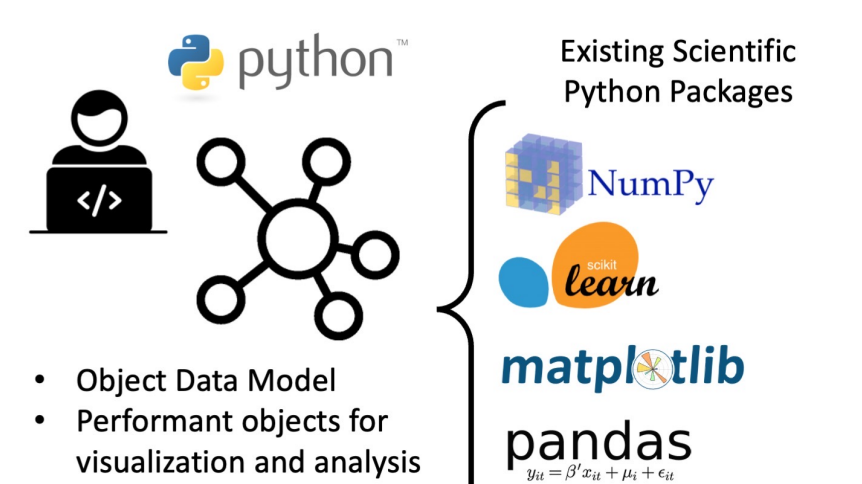
- Python mirror of the R dataRetrieval tool
 - Currently has most of the same functions
 - Very similar results
 - Collaborating with Timothy Hodson at USGS
- Example Jupyter Notebooks:
<https://www.hydroshare.org/resource/c97c32ecf59b4df90ef013030c54264/>
- <https://github.com/USGS-python/dataretrieval>



An advanced object data model that maps common water-related data types to performant data structures within the object-oriented Python language and analytical environment

A flexible water-data science object data model (hsmodels)

- Extending the HydroShare Resource Data Model to Python analysis environments
- Maps common water-related data types (HydroShare content types) to performant data structures within Python
- Load and stage data for visualization/analysis using common Python tools (pandas, matplotlib, etc.)
- <https://github.com/hydroshare/hsmodels>



HYDROLEARN

Educational module implementation in HydroLearn: Numbered steps indicate the workflow and location of essential module elements: (1) course landing page and links to course outline, (2) learning objectives, (3) module narrative, (4) code examples as interactive notebooks in the CUAHSI JupyterHub linked from HydroLearn, and (5) the technical assignment and associated rubric.