

HydroShare tools and recommended practices for sharing and publishing data and models in support of collaborative reproducible research

David Tarboton¹, Ray Idaszak², Jeffery Horsburgh¹, Daniel Ames³, Jonathan Goodall⁴, Alva Couch⁵, Pabitra Dash¹, Hong Yi², Christina Bandaragoda⁶, Anthony Castronova⁷, Martyn Clark⁸, and Shaowen Wang⁹

¹Utah State University

²Renaissance Computing Institute

³Brigham Young University

⁴University of Virginia Main Campus

⁵Tufts University

⁶University of Washington

⁷Consortium of Universities for the Advancement of Hydrological Science

⁸NCAR

⁹University of Illinois at Urbana Champaign

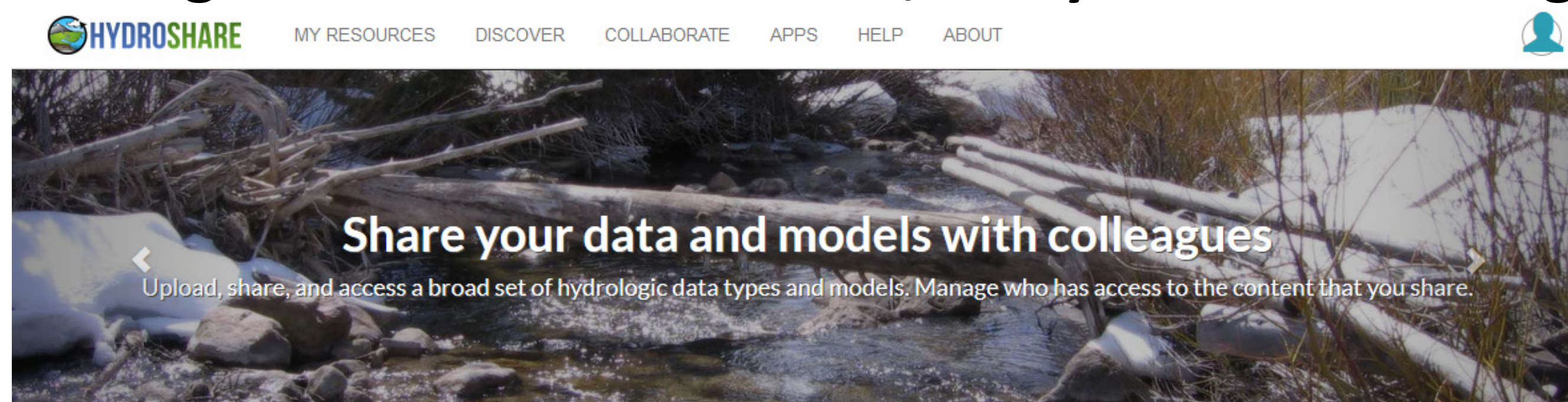
November 21, 2022

Abstract

HydroShare is a domain specific data and model repository operated by the Consortium of Universities for the Advancement of Hydrologic Science Inc. (CUAHSI) to advance hydrologic science by enabling individual researchers to more easily share products resulting from their research. The community platform supports, not just the scientific publication summarizing a study, but also the data, models and workflow scripts used to create the scientific publication and reproduce the results therein. HydroShare accepts data from anybody, and supports Findable, Accessible, Interoperable and Reusable (FAIR) principles. HydroShare is comprised of two sets of functionality: (1) a repository for users to share and publish data and models, collectively referred to as resources, in a variety of formats, and (2) tools (web apps) that can act on content in HydroShare and support web based access to compute capability. Together these serve as a platform for collaboration and computation that integrates data storage, organization, discovery, and analysis through web applications (web apps) and that allows researchers to employ services beyond the desktop to make data storage and manipulation more reliable and scalable, while improving their ability to collaborate and reproduce results. This presentation will describe the capabilities developed for HydroShare to support the full research data management life cycle. Data can be entered into HydroShare as soon as it is collected, and initially shared only with the team directly working on the data. As analysis proceeds, tools, scripts and models that act on the data to produce research results may be stored in HydroShare resources alongside the data. At the time of publication these resources may be permanently published and receive digital object identifiers and cited in research papers. Resources may themselves include citations to the research papers, thereby linking the publications to the supporting data, scripts and models. HydroShare design choices and capabilities for establishing relationships and versioning, based on simplicity, and ease of use, and some of the challenges encountered, will be discussed.

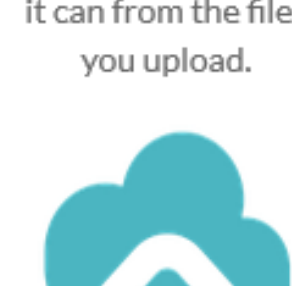
What is HydroShare ?

An online hydrologic information system for sharing data, models and code to enable more rapid advances in hydrologic understanding through collaborative research, analysis and modeling.



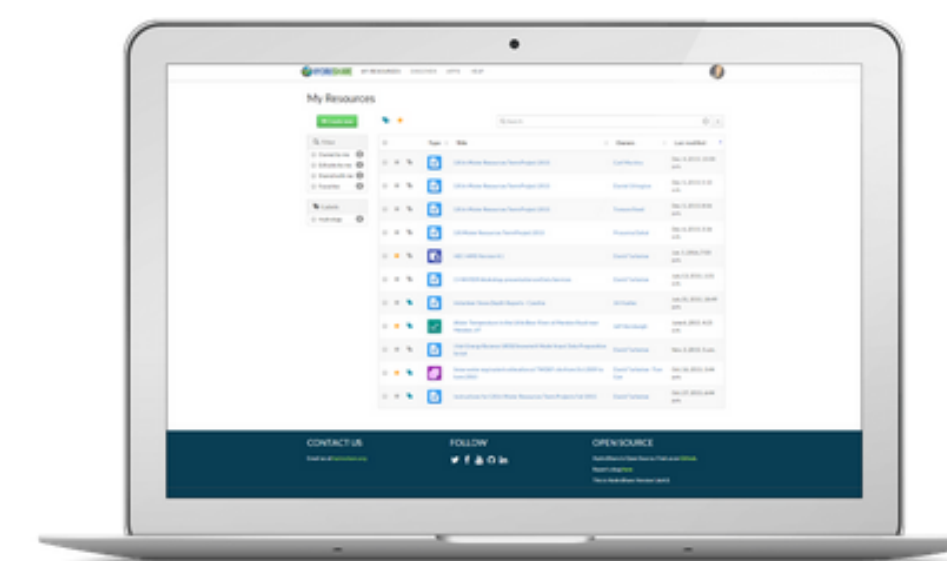
How it works

- Create data**
Collect your data using the same methods you use now. HydroShare supports a broad set of hydrologic data types.
- Upload to HydroShare**
Upload your data files to HydroShare through the web user interface. HydroShare will automatically extract as much metadata as it can from the files you upload.
- Describe with metadata**
Use HydroShare's simple metadata entry forms to finish describing your data so that your colleagues can find, access, and interpret it.
- Share with colleagues**
Choose who has access to the data and models you have uploaded to HydroShare. You can share with individual users or make your resources public for everyone to access.
- Permanently Publish**
Obtain a Digital Object Identifier (DOI) so your work can be easily cited. Reference related journal publications in your metadata.



What you can do with HydroShare

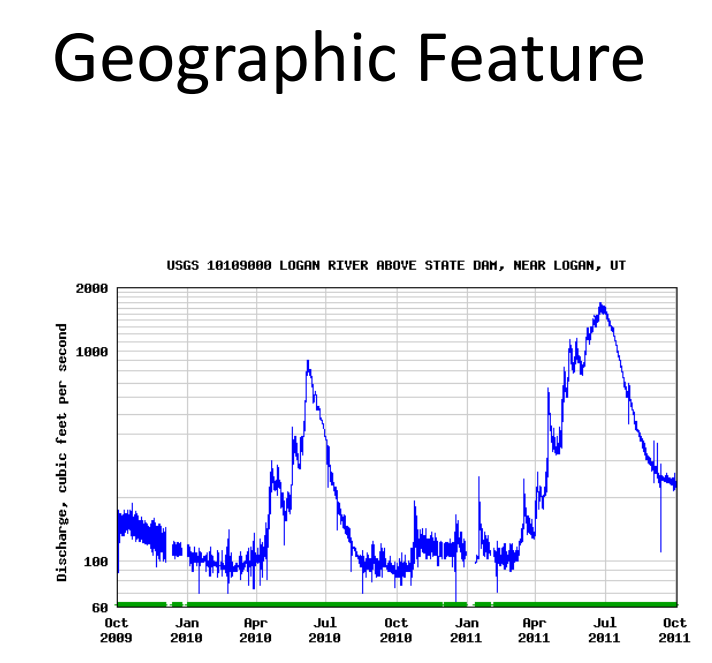
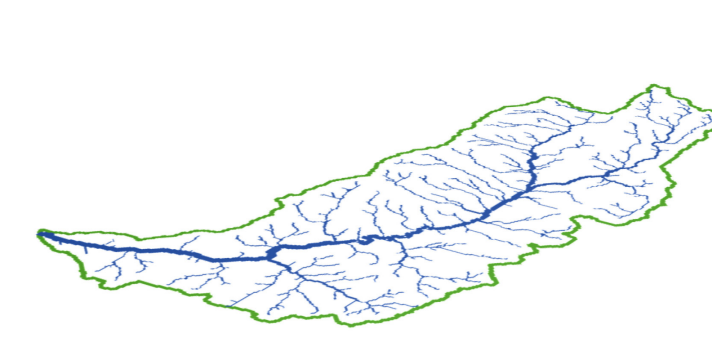
- ✓ Share your data and models with colleagues
- ✓ Manage who has access to the content that you share
- ✓ Share, access, visualize and manipulate a broad set of hydrologic data types and models
- ✓ Use the web services API to program automated and client access
- ✓ Publish data and models to meet the requirements of your data management plan
- ✓ Discover and access data and models published by others
- ✓ Use web apps to visualize, analyze and run models on data in HydroShare



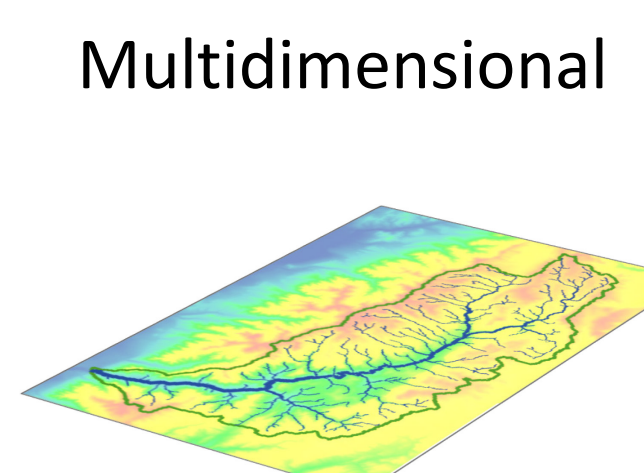
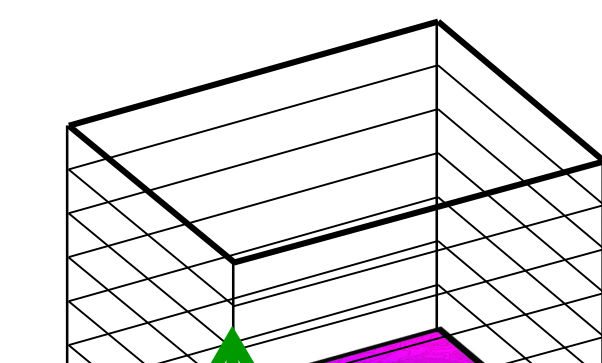
www.hydroshare.org

What can you store in HydroShare ?

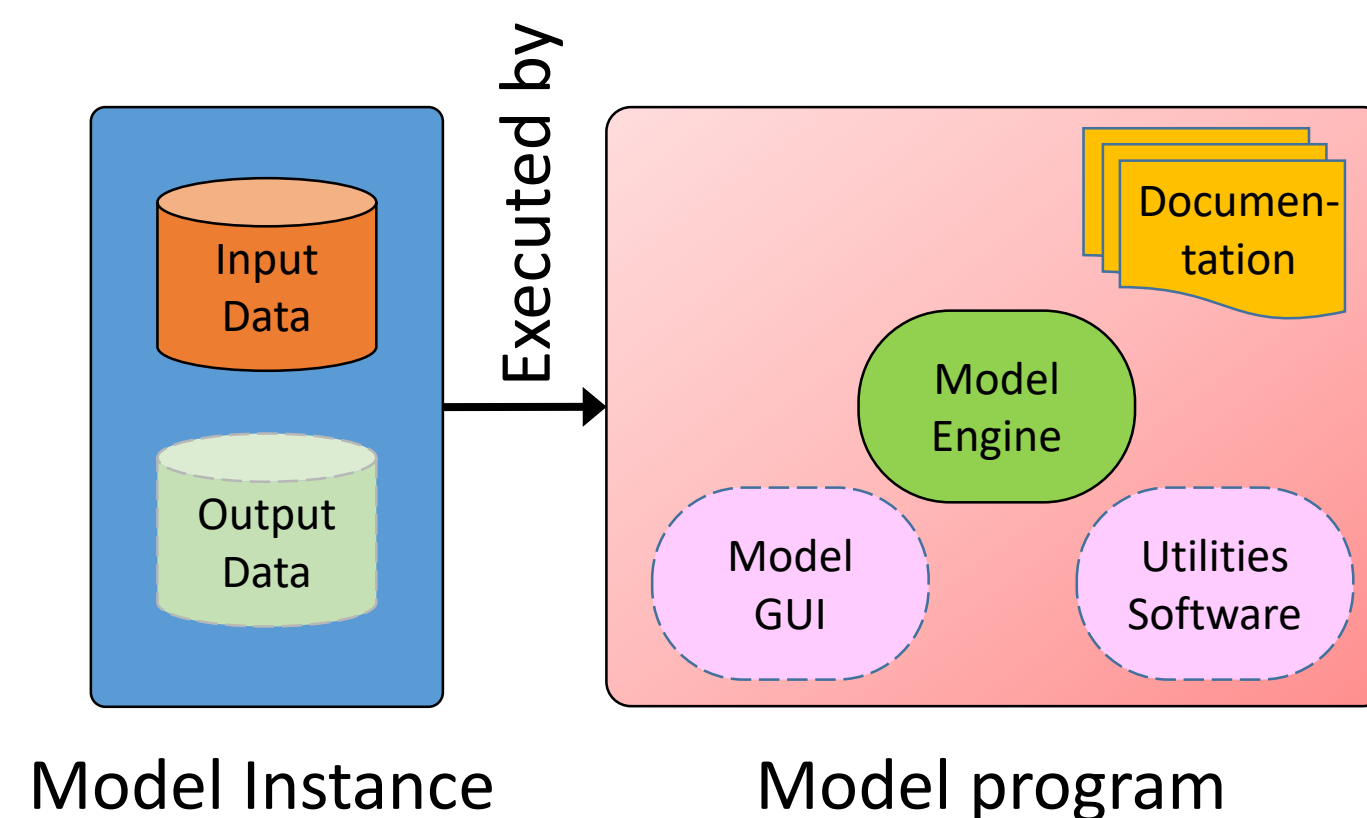
- In HydroShare, data and model files are stored as **resources**.
- HydroShare supports any file, including several specific data formats.
- Content “aggregations” hold data formats common in hydrology and support description with additional content specific metadata. Apps can act on specific content types.
- Collections group together multiple resources related to a project or study.
- Model Programs and Model Instances hold specific hydrologic models and associated data for application at a location.



Time Series



Geographic Raster



Why HydroShare ?

Collaboration: Share your data and model files; integrate information from multiple sources; organize individual, team, and group work.

Reproducibility, transparency and trust: Publish your work in any format, including data and models with a citable digital object identifier (DOI).

Do Science: Run Apps and models from a browser without installing software; access computational services for your big data and model analysis.

Learning: Use a platform where all students have access to the same functionality regardless of their computer.

HydroShare is a system to advance hydrologic science by enabling the community to more easily and freely share products resulting from their research, not just the scientific publication summarizing a study, but also the data and models used to create the scientific publication.

- Findable
- Accessible
- Interoperable
- Reusable

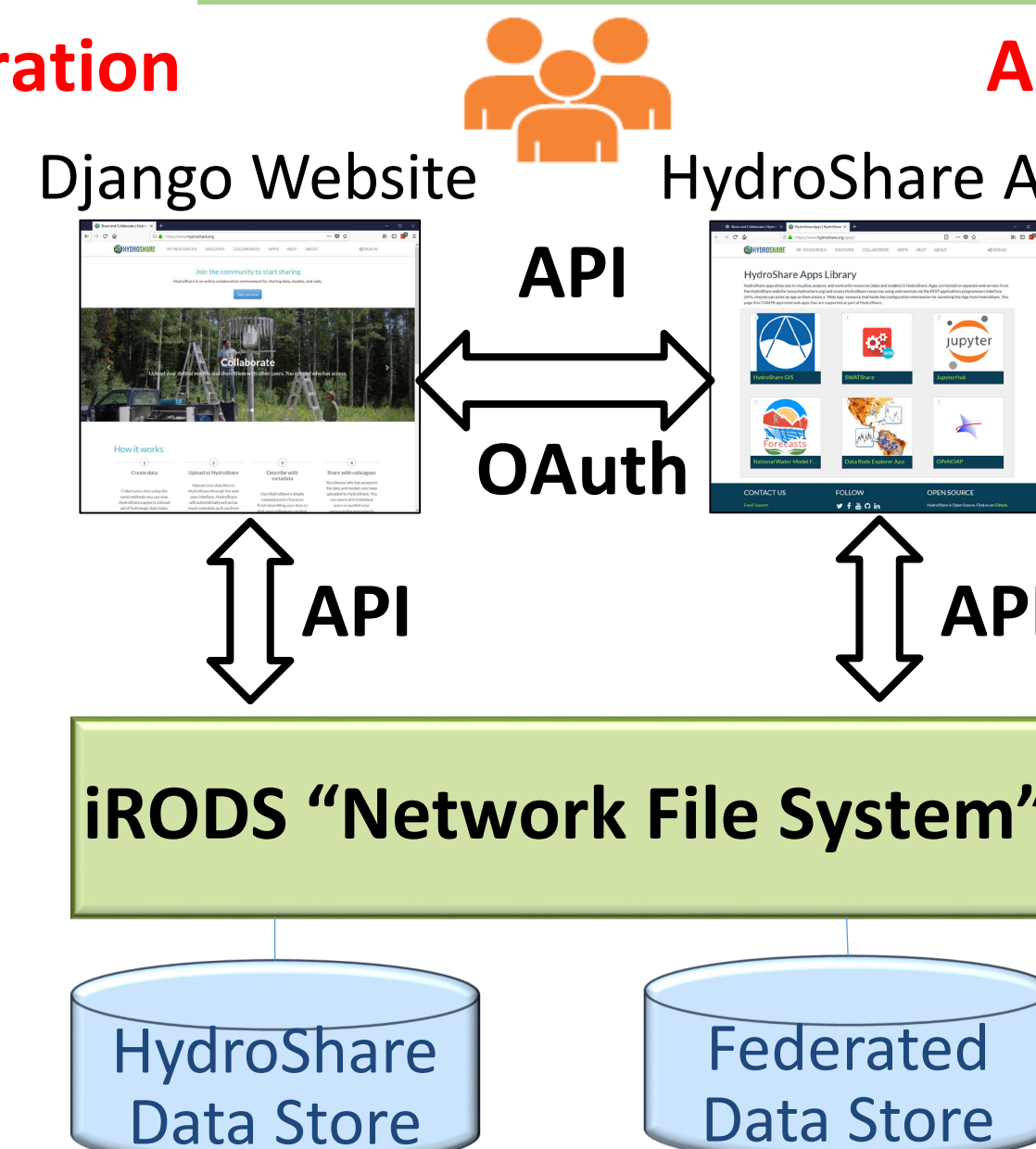


Design

Resource exploration

- Organize and annotate your content
- Manage access

Distributed file storage



Actions on Resources

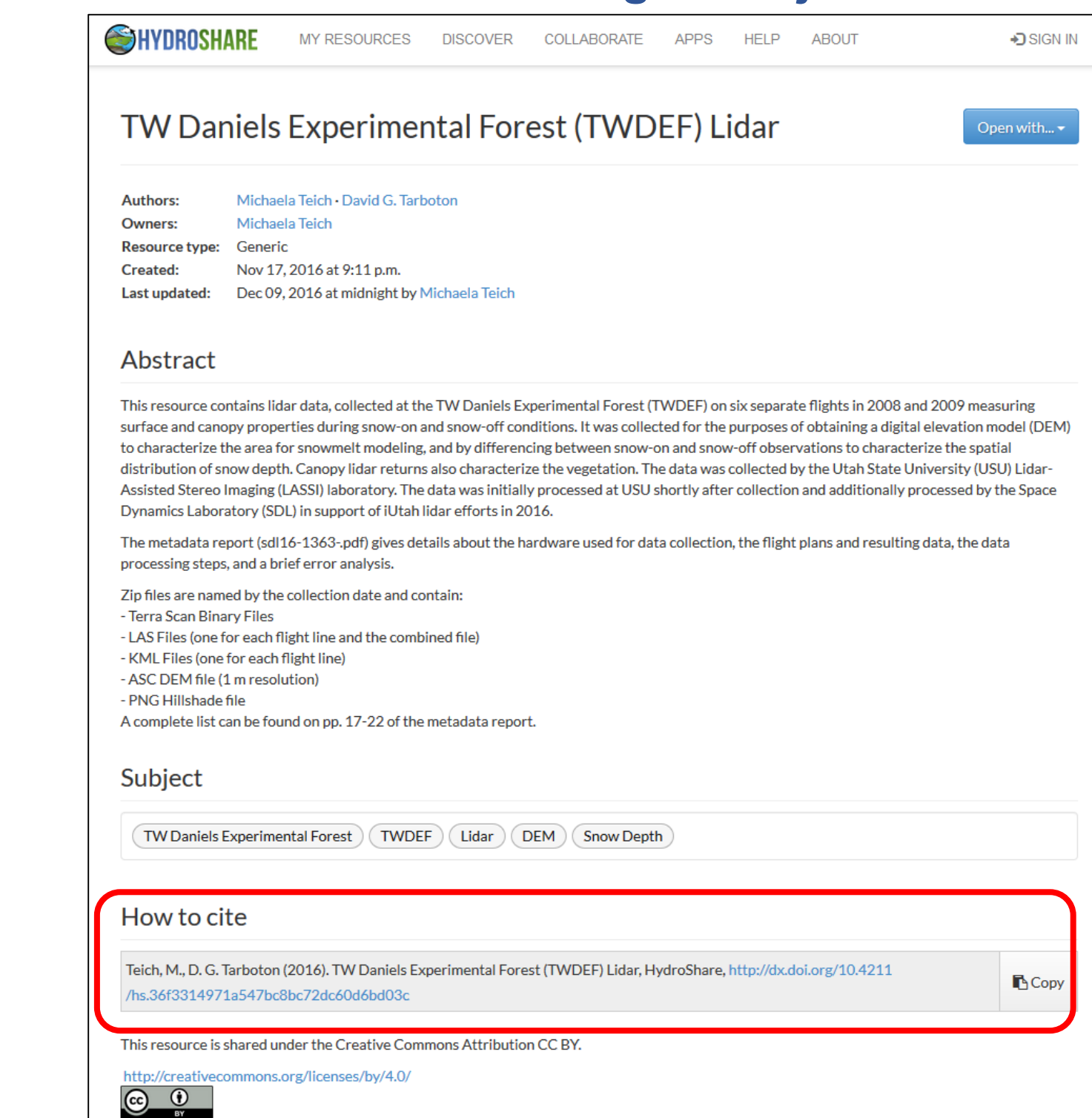
- Web software to operate on content you have access to (Apps)
- Extensibility: Anyone can set up a server/app platform (software service) to operate on HydroShare resources through iRODS and API

e.g. NCSA, TACC, USU

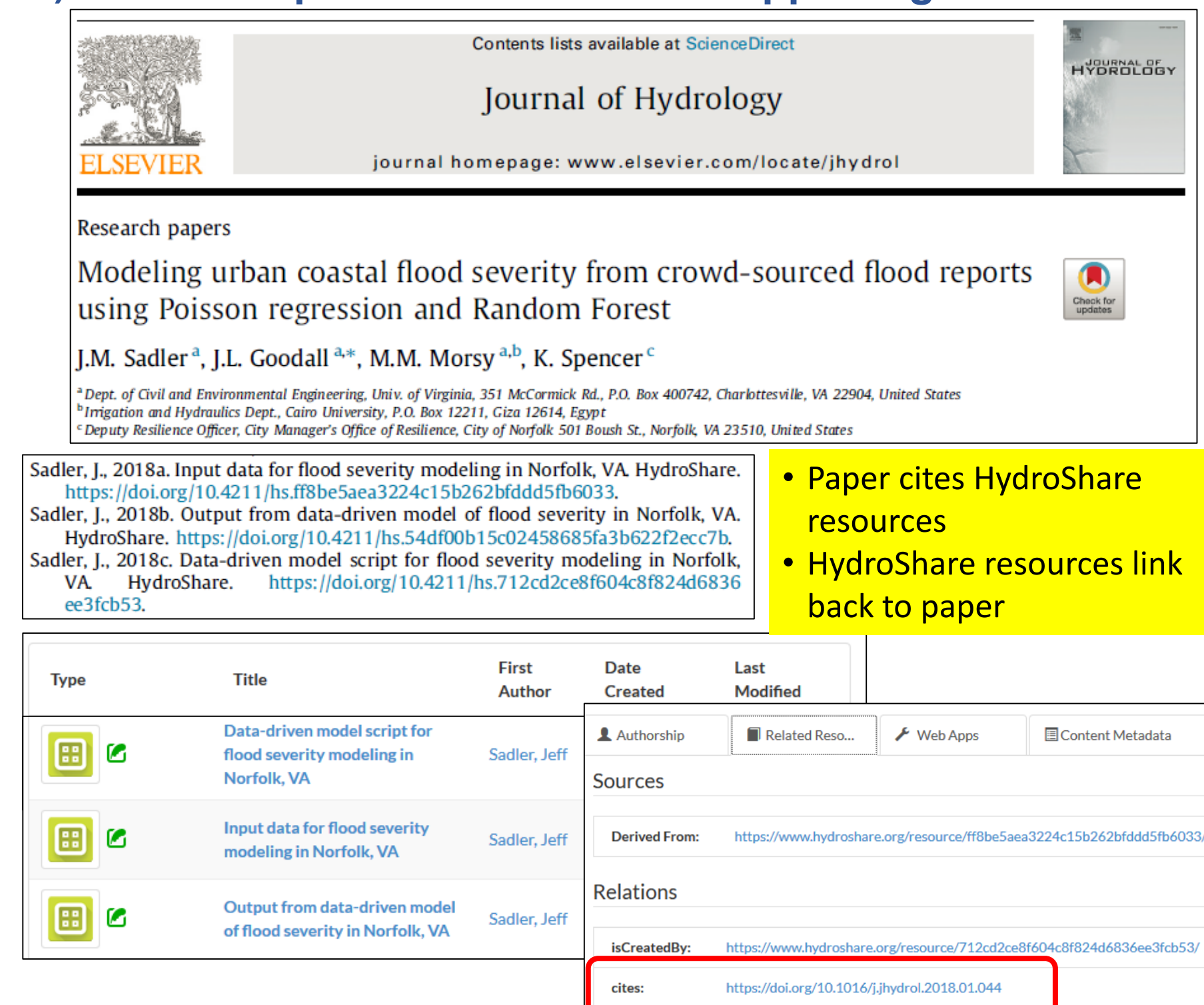
Moving towards fully web based hydrologic innovation environment

Publishing data and models

Publication with Citable Digital Object Identifier (DOI)

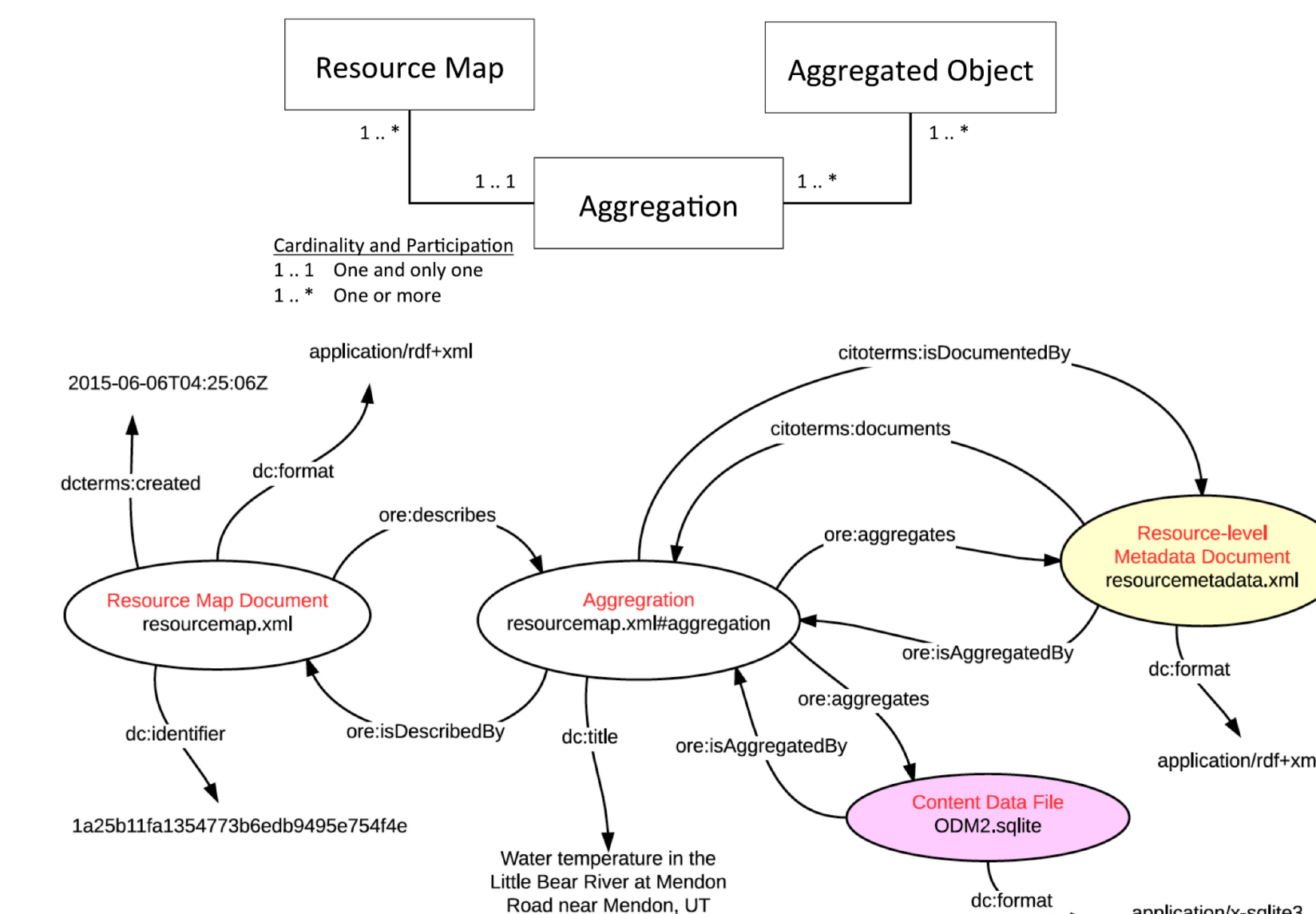


Link publications to their supporting data



- Paper cites HydroShare resources
- HydroShare resources link back to paper

OAI-ORE standard based Resource Data Model



Horsburgh, J. S., et al., (2016), "Hydroshare: Sharing Diverse Environmental Data Types and Models as Social Objects with Application to the Hydrology Domain," JAWRA, <http://dx.doi.org/10.1111/1752-1688.12363>.

Dublin Core machine readable metadata and data model to make data in HydroShare, Findable, Accessible, Interoperable, Reusable

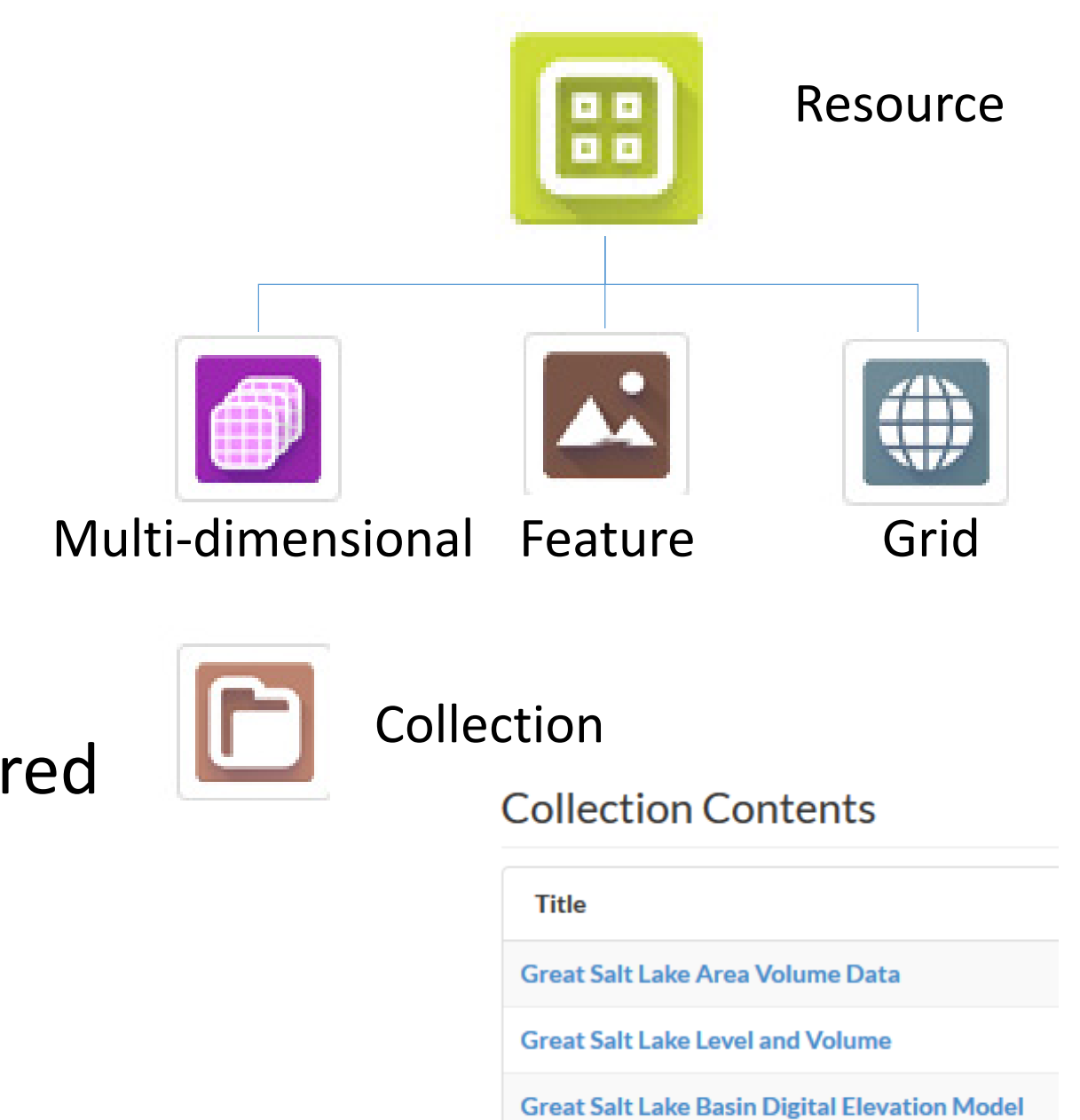


Resources, comprised of data and models, are framed as social objects, the basis for collaboration and interaction

Resource Organization

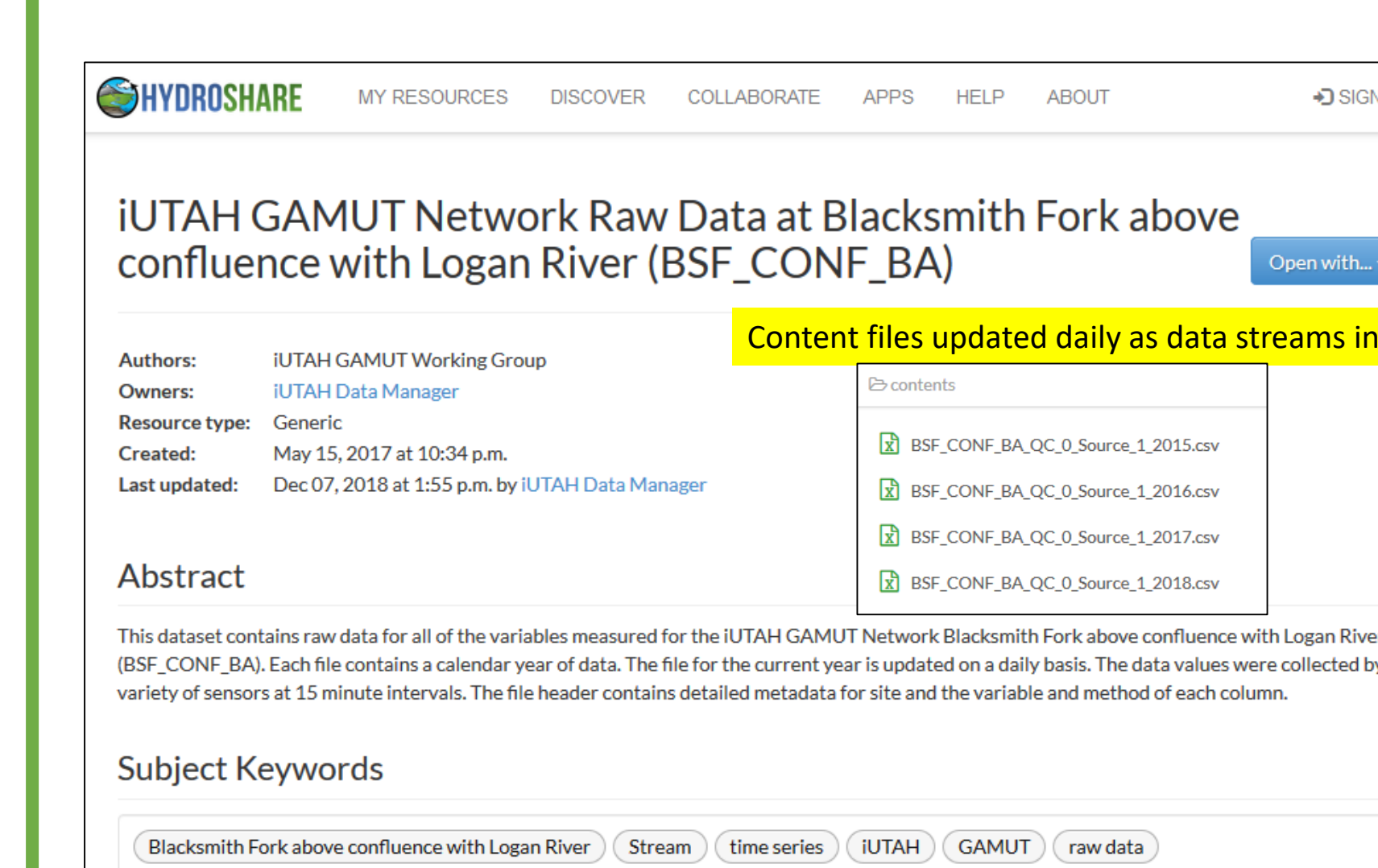
- A **resource** can hold multiple aggregations
 - Each being a different type of data with its own set of metadata
 - Managed as one discoverable resource
 - One set of access controls (Owners, Editors etc.)
 - One unique identifier
 - One set of resource level metadata

- A **collection** can hold multiple resources
 - Collections and their members may each be discovered separately
 - Unique keyword tags form informal collections (e.g. "AGU2018")

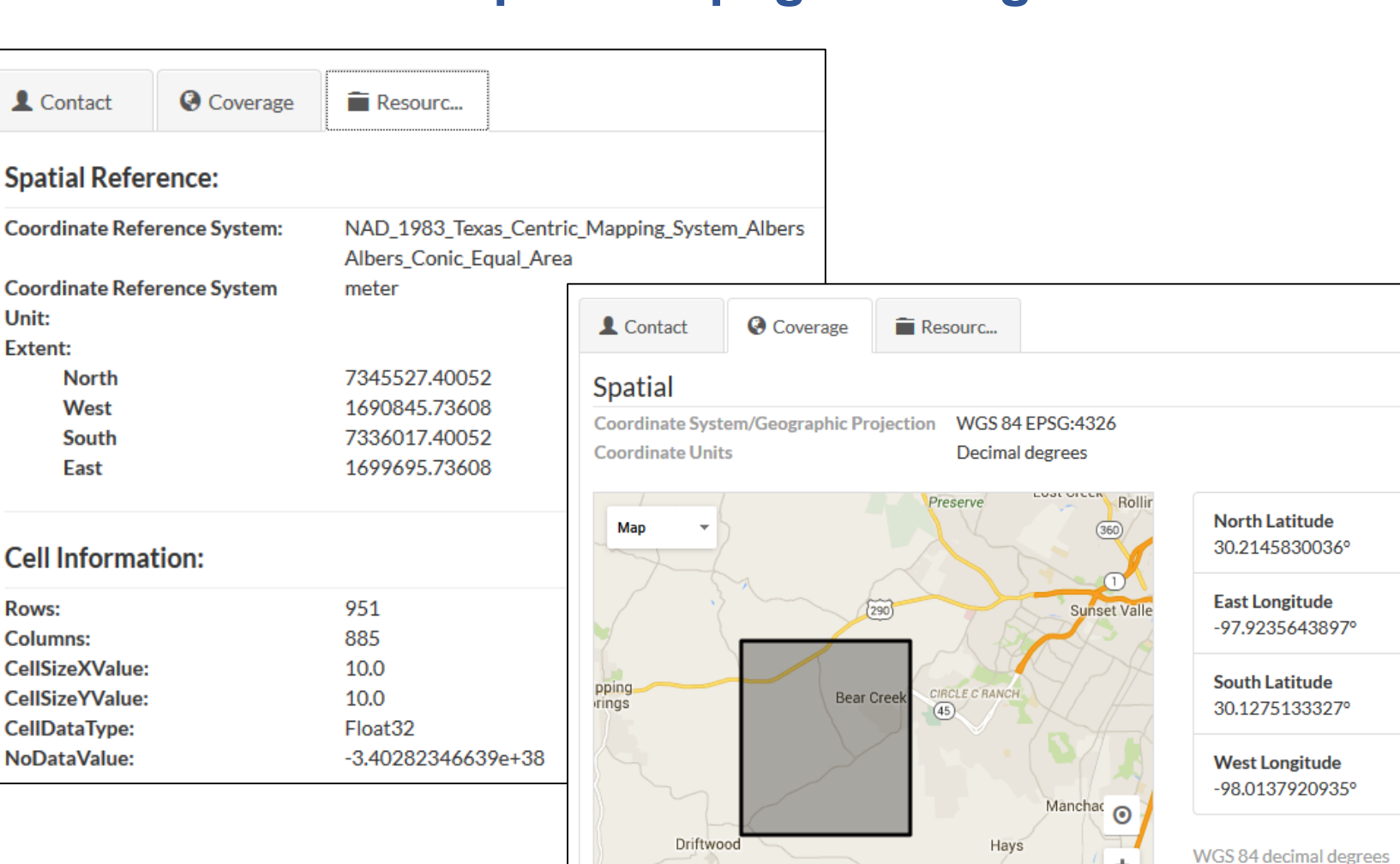


Key Functionality

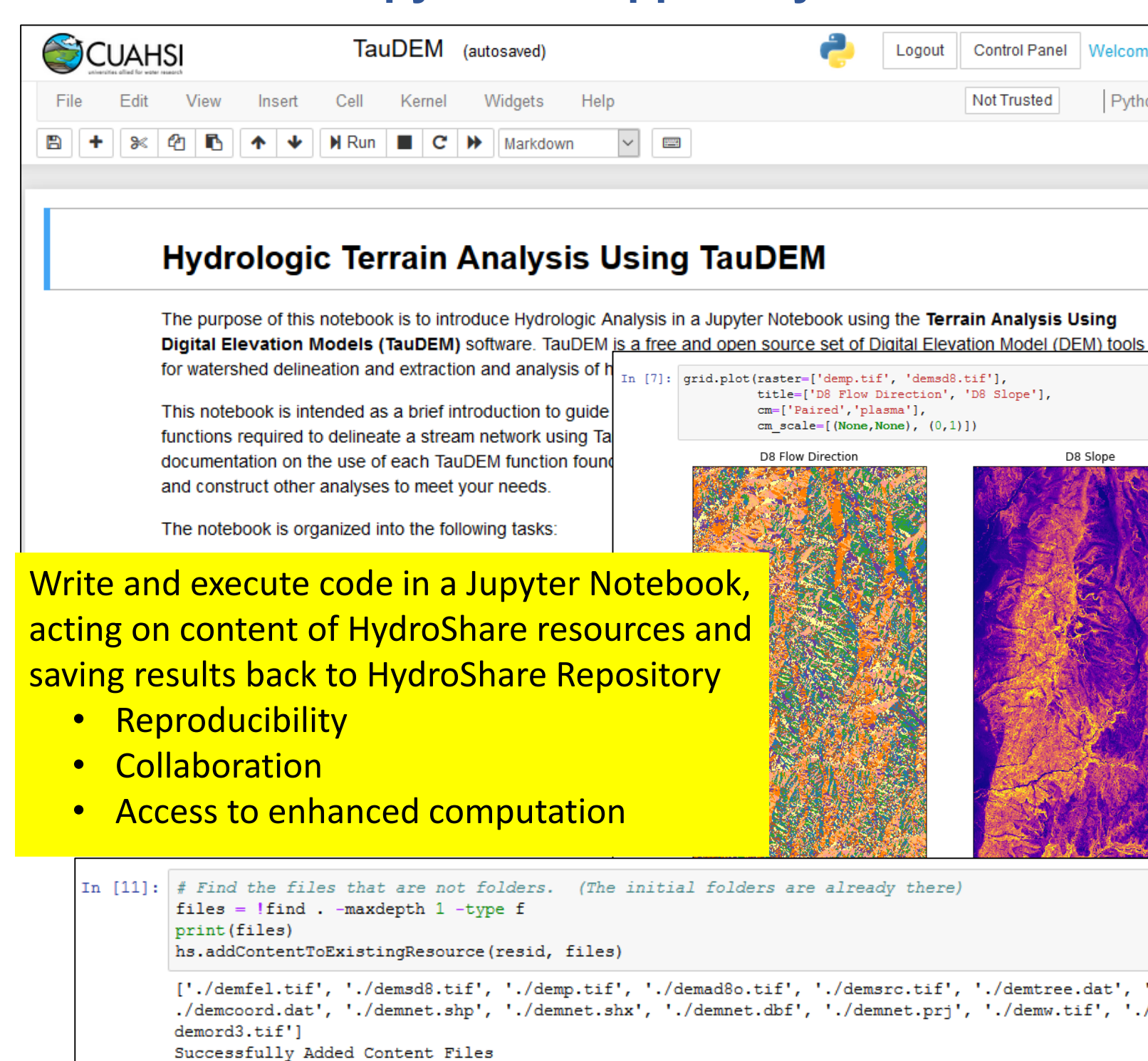
Data streamed into HydroShare as soon as it is collected



Metadata harvested automatically or captured via simple web page editing



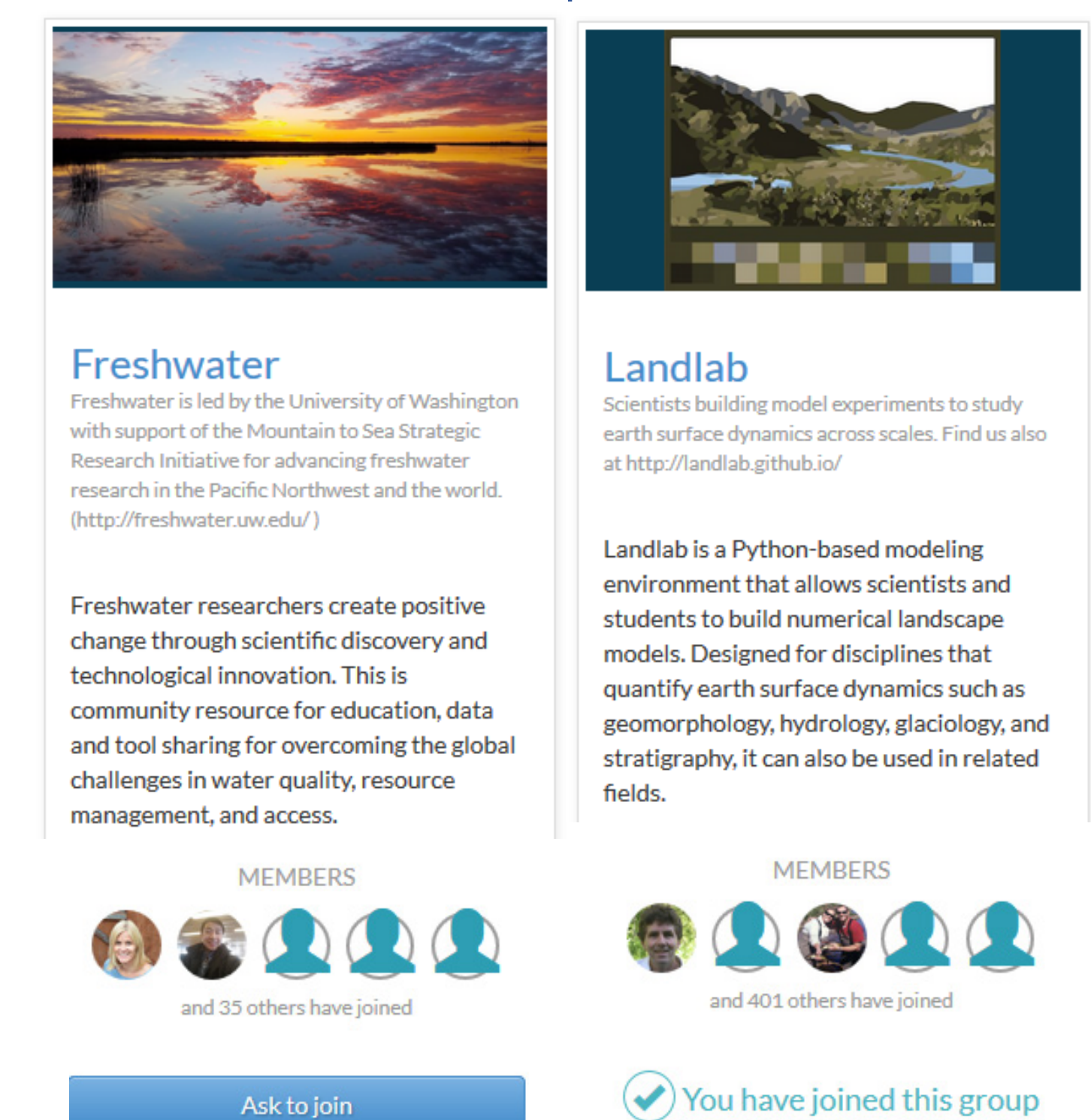
JupyterHub App Analysis



Write and execute code in a Jupyter Notebook, acting on content of HydroShare resources and saving results back to HydroShare Repository

- Reproducibility
- Collaboration
- Access to enhanced computation

Groups



OAC-1664061
OAC-1664018
OAC-1664119

HydroShare is operated by CUAHSI with ongoing development through a collaborative project among Utah State University, RENCI University of North Carolina, Brigham Young University, CyberGIS Center University of Illinois, Tufts, University of Virginia, NCAR, and University of Washington.



Universities Allied for Water Research

