

The I-GUIDE Cyberinfrastructure Platform: Supporting Open, Reproducible Convergence Science at Scale

Jeffery S. Horsburgh

Utah Water Research Laboratory

Utah State University

and

The I-GUIDE Core Cyberinfrastructure Capabilities and Services Team





Support: 2118329

Core CI Capabilities & Services Team



Anand Padmanabhan University of Illinois Urbana-Champaigr Carol (Xiaohui) Song





the Advancement of Hydrologic

Science, Inc. (CUAHSI)







Raiesh Kalvanam



I-GUIDE

Furgan Baig University of Illinois Urbana-Champaign

Anthony Castronova The Consortium of Universities for Irene Garousi-Neiad CUAHSI

Jeff Horsburgh Utah State University

Purdue University

Erick Li University of Illinois Urbana-Champaign













Noah Samuel Oller Smit

Objective: Integrate distributed geospatial data capabilities and advanced CI to form a composable and open I-GUIDE platform to accelerate scientific workflows and support education and workforce development as well as broader community engagement

Josh Lieberman The Open Geospatial Consortium **Alexander Michels**

University of Illinois Urbana-

Champaign

Steven Manson

Mohan Ramamurthy University of Minnesota Twin Cities University Corporation for Atmospheric Research (UCAR)

Ingo Simonis

The Open Geospatial Consortium Purdue University



Sina Taghavikish



The Open Geospatial Consortium

Purdue University

What is the I-GUIDE "Platform"?

- Existing distributed datasets and community cyberinfrastructure
- We want to get to convergence research
- The platform is the bridging infrastructure between the two



Confusing terminology...

- The "buzz words"
 - I-GUIDE Platform?
 - I-GUIDE Gateway?
 - Geospatial dataon-demand?

• Confusing – how are these things different than the "Platform"?



Some challenges in communication

- Convergence Science Catalyst (CSC) teams want to use the cyberinfrastructure, but they don't know what it is, how to use it, or how it improves their life
- The Core Cyberinfrastructure Capabilities and Services Team wants to build effective cyberinfrastructure, but we don't fully understand the science use cases

We need some common ground to move forward together!



The purpose of this VCO

- Describe the I-GUIDE Platform in a bit more detail
- Provide our vision and approach for the Platform
- Describe the principles that are guiding our development
- Describe some of the challenges we are working to address

I-GUIDE Platform Vision

From the perspective of the I-GUIDE cyberinfrastructure team, we want to:

Support open, reproducible convergence science <u>at scale</u> using distributed geospatial datasets that may be too large to store in a single location. The I-GUIDE Platform aims to help reduce the time to science and enable more open, accessible, and reproducible computational workflows.

"At scale"

- Geospatial convergence science faces multiple scale-related challenges:
 - <u>Data</u>: How to efficiently store, manage, and stage data for modeling and analysis? How to share results?
 - <u>**Computation**</u>: How to get access to adequate computational resources?
 - <u>**Connection**</u>: How to efficiently connect data and computation?
 - <u>**Reproducibility</u>**: How to make all of these accessible for reproducing computational work?</u>



More services are now cloud based Composing workflows that bridge clouds is hard Advancing Innovation amazon S3 **≊USGS** Compute Storage Data \square **HYDROSHARE** jupyter Analysis Visualization Repository

Our approach to building the Platform

- Use existing cyberinfrastructure components where we can
- Find ways to facilitate "composing" workflows that use those components
- Build missing pieces where necessary
- Explore what we can do by engaging with use cases coming out of the convergence science teams

Guiding principles for working together

- We need some principles to guide our cyberinfrastructure development
- To find the common ground, these also need to be shared by I-GUIDE's convergence science teams!
- 1. I-GUIDE convergence science should be open meaning shared, transparent, and accessible.
- 2. Inputs and outputs of analyses along with workflows used to generate them should be shared in open repositories or openly accessible storage locations.
- 3. Computational workflows should be defined in a way that they can be documented, shared, and reproduced/repeated using accessible computational resources.

Applying these principles will likely require some changes to the way we work.

What is driving these principles?



- Data (broadly defined) are primary research products!
- Findable: Data have sufficient metadata and a unique, persistent identifier making data discoverable on the Web
- <u>Accessible</u>: Metadata and data are understandable to humans and machines and are available via a trusted repository
- Interoperable: Metadata use formal community standards
- <u>Reusable</u>: Data have clear metadata, usage license, and information about provenance

The extent to which data are FAIR affects their value and extent of reuse!

Wilkinson, M. D. et al. (2016). The FAIR Guiding Principles for scientific data management and stewardship. Scientific Data, 3:160018, https://doi.org/10.1038/sdata.2016.18.

Funder requirements

NSF

NSF's general Data Sharing Policy¹:

"NSF-funded investigators are expected to share with other researchers, at no more than incremental cost and within a reasonable time, the primary data, samples, physical collections and other supporting materials created or gathered in the course of work under NSF awards."

• From NSF's Division of Earth Sciences Data and Sample Policy²:

"Possible types of "data" to be addressed in the DMP include, but are not limited to: observational, experimental, analytical, and model outputs; derived and compiled datasets; software and code; educational materials; and any other relevant digital products resulting from the project."

¹ <u>https://new.nsf.gov/funding/data-management-plan#nsfs-data-sharing-policy-1c8</u> ² <u>https://www.psf.gov/funding/data-maliging/gata-fata-sharing-policy-1c8</u>

² https://www.nsf.gov/geo/geo-data-policies/ear/ear-data-policy-jul2023.pdf

Publisher requirements



"AGU requires that the underlying data and/or software needed to understand, evaluate, and build upon the reported research be available at the time of peer review and publication. Additionally, authors should make available software that has a significant impact on the research. This entails:

- 1. Depositing the data and software in a community accepted, trusted repository, as appropriate, and preferably with a DOI
- 2. Including an Availability Statement as a separate paragraph in the Open Research section explaining to the reader where and how to access the data and software
- 3. And including citation(s) to the deposited data and software, in the Reference Section¹."

¹ <u>https://www.agu.org/Publish-with-AGU/Publish/Author-Resources/Data-and-Software-for-Authors</u>

Publisher incentives

To incentivize authors to make their results more reproducible, the Journal will **publish technical papers** and case studies with verified reproducible results **open access free** to the authors for the next year."



access either free or for a reduced fee, as funds are available. The *Journal* will also recognize papers with reproducible results in a new special collection and offer two new annual reproducibility awards for authors and the people who assess the reproducibility of results.

Add a Reproducible Results Section

Authors of any article type may add an optional "Reproducible Results" section to their manuscript if all research materials listed in the "Data Availability Statement" section are publicly available in a repository or online. Next, authors must follow best practices to improve result reproducibility (Rosenberg et al. 2020). One best practice is to have a colleague, student, or other person not affiliated with the study reproduce manuscript results. In the "Reproducible Results" section, authors must list which figures, tables, or other results in their manuscript have been reproduced and by whom. This demonstrates that the authors have identified and addressed any bugs, missing materials, unclear directions, and other glitches that often befall initial versions of repositories and prevent reproduction. A "Reproducible Results" section for an initial

PDF

From the I-GUIDE Data Management Plan

"The institute itself <u>will strictly adhere to NSF rules</u>, and share with other researchers, at no additional cost and within a reasonable time, the primary data, software and other related materials created or gathered during the project execution."

"We will archive digital data used in our experiments, as well as the associated software and programs generated by the institute and make such materials available to pertinent communities. For example, Hydroshare, will be used for hosting hydrology related datasets and associated models."

"From the point of view of this institute, such dissemination of data is crucial to stimulate new advances as quickly as possible and to allow for prompt evaluation of the results by the pertinent communities."



What do we mean by "open"?

• An example workflow for a dataset product:



and uploads to reputable repository

iterates collaboratively

published and assigned DOI

dataset in paper

Result: Research artifacts (data and analyses) are openly shared, accessible, and citable.

Open repositories

- Choose a repository that provides:
 - Data registration/publication with a persistent, globally-unique identifier such as a digital object identifier (DOI)
 - Free access to the data
 - A landing page that provides metadata
 - Support for versioning













Data



- Web-based repository and Hydrologic Information System
- Operated by the Consortium of Universities for the Advancement of Hydrologic Science, Inc. (CUAHSI)
- Creating and sharing data and models using a variety of file formats and flexible metadata
- Formal publication of data and models with linkages to literature



From our prior VCO about Hydroshare:

HYDROSHARE enables transparency









User creates dataset and uploads to Hydroshare Research team iterates collaboratively

Final version is published and assigned DOI User cites dataset in paper

What happens when this breaks down?



Data constraints

- Limitations on data size represent a constraint imposed by current repository technology
- Some important datasets are too large to store in existing repositories like HydroShare
- Some models generate results that overwhelm repository storage
- As we scale up, we need other options!

Despite policies, journals don't have all the answers either!

When it comes to large datasets, we are often asked by authors and editors how they should preserve the data. These questions come via datahelp@agu.org and our data and software guidance discussions. Spoilers, there are no easy answers, yet! Here we offer our experience, share the current limitations, and the approaches we recommend with what is possible right now.

AGU requires that primary and processed data used for your research should be preserved and made available. This can range from observational data to the data used to generate your figures. The raw data may be needed, but usually, the processed or refined data that support and lead to the described results and allow other readers to assess your conclusions and build off your work should be preserved.

For data that is large, over 1 Terabyte (TB), authors run into the challenge of finding a suitable repository. Many repositories have file size limitations but also costs associated with deposits over certain limits. This generalist repository comparison chart provides an overview of the limitations. Discipline-specific and institutional repositories are often a place to turn to for assistance with preserving large data but they also have limitations and potential costs. This emphasizes the importance of avoiding surprises at the time of publication by:

- planning ahead (e.g. type of data, size of files),
- finding a suitable repository (e.g. discipline, institution, general),
- knowing the limitations of the repository (e.g. upload size), and
- determining costs ahead of time (e.g. size, years of preservation).

In other cases where the data is too large and complex to move and deposit, for instance, the simulation data from models and associated workflows running at computing clusters, the discussion turns to what data should be preserved. In AGU'S Data and Software Guidance for Authors we outline the decisions an author needs to make (see Models and Simulations) and we go into further detail in our journal specific guidance as well. There is also a group working on a Rubric for Models and Model Data — Best Practices for Preservation and Replicability.

Some research computing facilities and institutions provide sharing and access mechanisms for such

What does it mean to be reproducible?

Essawy, B., J. Goodall, D. Voce, M. Morsy, J. Sadler, Y. D. Choi, D. Tarboton and T. Malik (2020), A taxonomy for reproducible and replicable research in environmental modeling, *Environmental Modelling & Software*, 134:104753, <u>https://doi.org/10.1016/j.envsoft.2020.104753</u>

From our prior VCO about Hydroshare: Cloud Computing Makes Resources Actionable

Sample Computational Resource					Open with	C File	File Edit View Run Kernel Tabs Settings Help Residence unvelopment Basic-ration-curve involutional settings Help Residence unvelopment Residence unvelopment						
				Dat	ta Series Viewer	0	+ ☆ □ □ ► ■ C [143]:	▹ Markdown ∨ agency_cd site_no tz_cd s	tage_ft 65347_	_00065_cd disc	charge_cfs 65349_	_00060_cd	眷 Python 3 (ipykernel)
Authors: Owners:	Anthony Michael Castronova Anthony M. Castronova	Sharing Status: Views:	Public 852	Jupyter Jup	pyterHub - CUAHSI DEV	≡ *	2020-06-04 00:00:00 2020-06-04 00:15:00 2020-06-04 00:30:00 2020-06-04 00:45:00	USGS 01011000 EDT USGS 01011000 EDT USGS 01011000 EDT USGS 01011000 EDT	3.12 3.12 3.12 3.12	A A A	1270.0 1270.0 1270.0 1270.0	A A A	
Type: Storage:	Resource The size of this resource is 2.4 MB	Downloads: +1 Votes:	102 Be the	🥠 ма	ATLAB Online		2020-06-04 01:00:00 2021-06-11 11:00:00	USGS 01011000 EDT USGS 01011000 EDT	3.12 2.58	A P	1270.0 712.0	A P	
Last updated: Citation:	Feb 22, 2021 at 5:29 p.m. Anthony M. Castronova See how to cite this resource	Comments:	NO CC	Cyt	berGIS-Jupyter for Water		2021-06-11 11:15:00 2021-06-11 11:30:00 2021-06-11 11:45:00 2021-06-11 12:00:00	USGS 01011000 EDT USGS 01011000 EDT USGS 01011000 EDT USGS 01011000 EDT	2.58 2.58 2.57 2.57	P P P	712.0 712.0 703.0 703.0	P P P	
Content	Q. Search current directory		0	TH	IREDDS Data Server		<pre>25932 rows x 7 column [144]: fig, ax = plt.subp p1 = df.stage_ft.p p2 = df.discharge_ av1 label = av[0]</pre>	s ots(2,1) ot(ax=ax[0], color='g') fs.plot(ax=ax[1], color=' et vlabel('Stage ft')	b')				
Contents	1ATLAB_Analysis.mlx		٩				ax2_label = ax[1]. fig.tight_layout()	<pre>ct_ylabel('Discharge, cfs</pre>	·,				
P R	ython_Analysis.ipynb _Analysis.R						25	2070-11 2021-03 2021-03 2021 datetime	55				
Csv C	bservation_data.csv eadme.md					Simple	<u>لاً بالم</u> <u>م</u> <u>م</u> <u>م</u> <u>م</u> <u>م</u> <u>م</u> <u>م</u> <u></u>	3 (ipykernel) Idle	hum .		Mode	e: Command 🛞	Ln 1, Col 1 basic-rating-cur

https://jupyterhub.cuahsi.org

What happens when this breaks down?

I want to create a Jupyter Notebook that enables people to run subsets of the National Water Model!

This workflow requires:

- Staging forcing data (~30 TB) and static data (~20 GB) somewhere
- Executing a domain subsetter program to get static and forcing data for a model domain
- Storing the subset of data for the model run (~5 GB)
- Executing the model on the resulting subset using a high-performance computational resource
- Storing the resulting model output data (~450 GB)
- Providing visualization and analysis of model outputs
- Rinse and repeat...

Computational constraints

- Availability and accessibility of computational resources represent constraints on the reproducibility of computational work
 - Scientists may not know which environment to use or how to use it
 - They may not have access to adequate computational environments
 - They may not be able to give others access to the environment they used
 - The computational environment they can access may not have convenient access to the data they need

How is the I-GUIDE Platform addressing these constraints?

- General I-GUIDE Platform Functionality
 - **Data staging and preparation**: Where can we put large datasets to make them available for processing and modeling?
 - **Data exploration**: How can we interact with large datasets given that we can't download (or don't want to move) the whole thing?
 - **Data integration/wrangling**: How can we subset large datasets? Aggregate in space or time? Combine with other datasets?
 - <u>Model/workflow execution</u>: Where and how can we run models over large domains or workflows that are computationally intensive?
 - **Data/workflow/results sharing**: What do we do with the results when we are done? How can we give others access?

I-GUIDE Platform Design

Two main user interfaces

I-GUIDE Catalog: Accessing shared data, workflows, etc.

I-GUIDE JupyterHub(s): Multiple options for data exploration, integration, wrangling, model/workflow execution

😑 🔵 💭 📿 JupyterHub	× +		
← → ଫ ଲ 💷 jupyter.igui	de.Illinois.edu/hub/spawn	* •	
C Jupyterhub Home Tok	en	jeff.horsburgh@usu.edu	(+ Logout
	Server Options		
	I-GUIDE Platform		
	Default Environment for the I-GUIDE Platform (most users should choose this)		
	I-GUIDE Esri Summer School 2023		
	I-GUIDE Esri Summer School 2023		
	Start		

Current Prototype: <u>https://iguide.cuahsi.io</u>

Example: https://jupyter.iguide.illinois.edu/

How do we envision people using the Platform?

Stage your data in an appropriate

storage location or repository

Choose an appropriate computational resource

Deposit data, notebooks, results in a repository for publication

Register your products with the I-GUIDE catalog

These steps aren't trivial - some iteration may be required

Computational Notebook

- Store and retrieve content
- Staging data for analysis
- Sharing data/notebooks/content
- Launching computational notebooks
- Permanent publication

- Moving data
- Staging data for modeling
- Matching storage with computation
- Large data output

Link to HPC where needed •

Launch "actionable products" (e.g., Notebooks)

Example Platform applications

 Jupyter Notebooks that demonstrate capabilities of the platform

Where do we need to improve the Platform?

- The entry point for using the Platform isn't clear where do I start?
- There are a lot of options
 - Which storage should I use?
 - Which JupyterHub should I use?
 - Which repository should I use?
- We developed a generalized system, but we need to further test and exercise with I-GUIDE use cases
- In other words do the components we have assembled meet the needs of the CSC teams?

Why invest in using the I-GUIDE Platform together?

- Greater transparency for our scientific work and increased trust in results
- Greater opportunity for enabling reproducibility of results
 - <u>Reproducibility by design</u> rather than requirement
- Published research products (with citable DOIs)
- Tools for sharing and collaboration
- Facilitated access to repositories and computational resources
- Tools for teaching and active learning
- Compliance with funding agency and journal publisher requirements
- Because we said we would!

Questions?

Jeffery S. Horsburgh jeff.horsburgh@usu.edu

Utah Water Research Laboratory UtahStateUniversity

Support: 2118329