



CIROH 1st Annual Training and Developers Conference 2023
Hands-On Demo of the CIROH Integrated Evaluation System Prototype Workshop
Wednesday, May 17th at 1:30 PM

Katie van Werkhoven
Matthew Denno
John Park
Chris Townsend



TEEHR

Tools for Exploratory Evaluation in Hydrologic Research

Background and Motivation

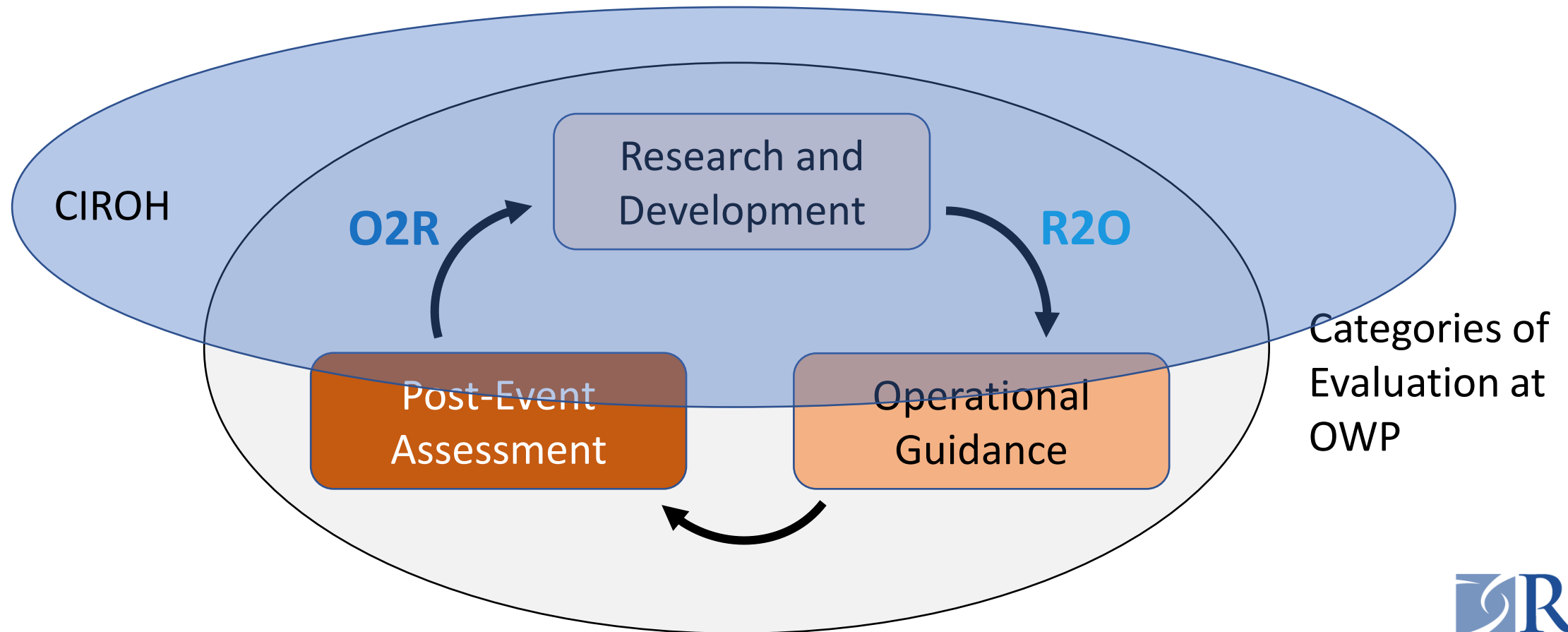


What is
Evaluation?



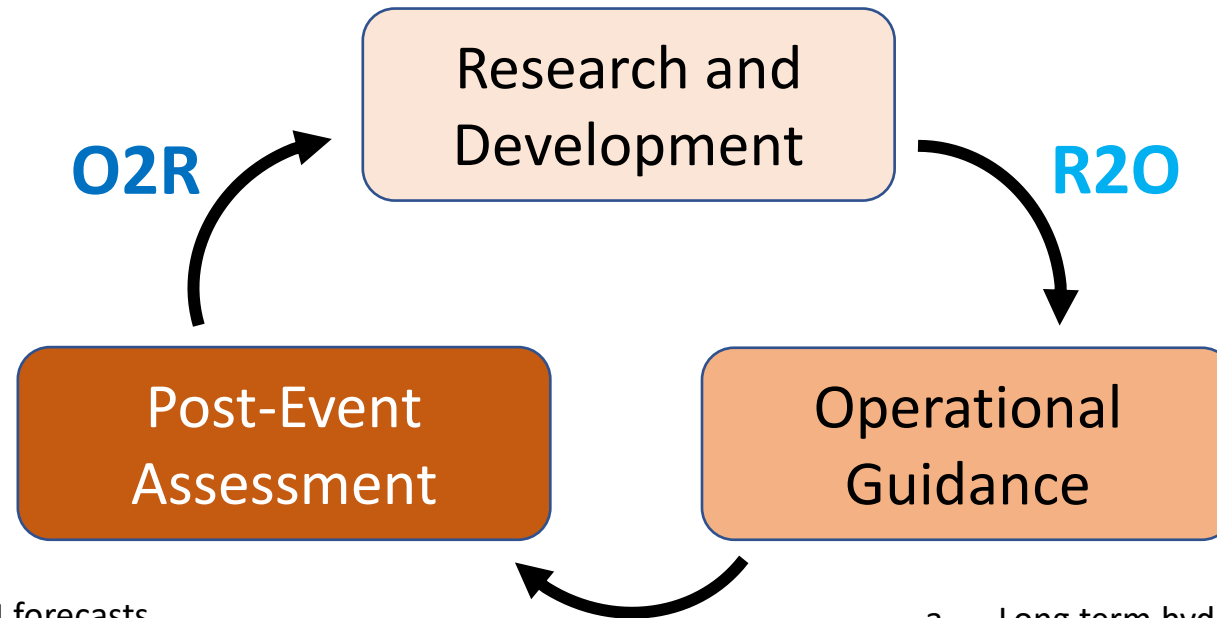
What types of evaluation are we talking about?

An assessment of the quality of hydrologic forecasts, or *component* of the forecast process, with respect to some baseline



A few specific examples

- a. Hydrologic model component development (snow, land surface, subsurface, reservoir, channel...)
- b. Flood inundation model/methods
- c. Operational forecasting methods/approaches (DA method, forecast forcing, ensemble methods...)

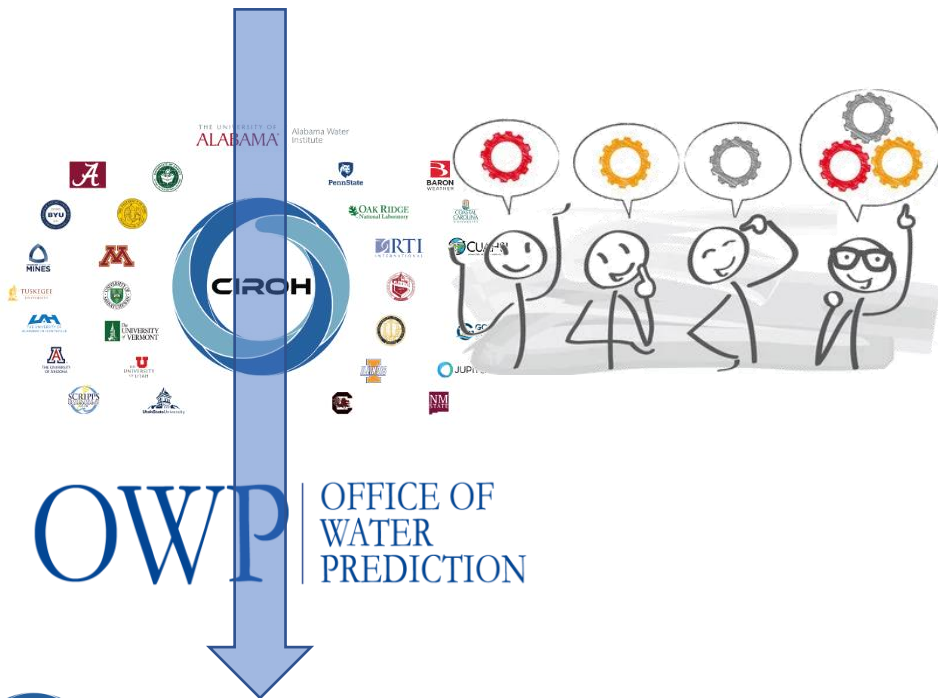


- a. Qualitative reviews of NWM forecasts
- b. Precip forecast performance assessment
- c. Hydro forecast performance assessment
- d. Hydro model performance assessment
- e. FIM performance assessment

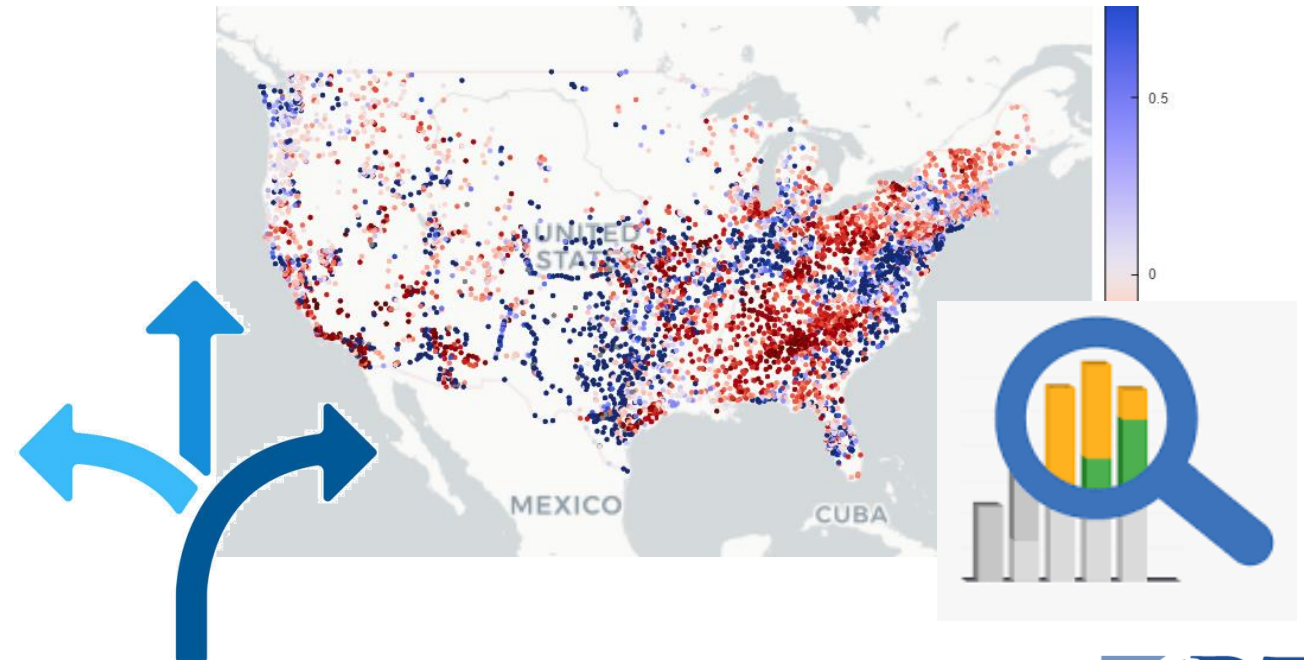
- a. Long term hydrologic model error trends
- b. Forecast forcing (QPF) trends
- c. Hydrologic reforecast/hindcast trends
- d. Current model condition assessment

Why a CIROH Evaluation System?

Foster Consistency and
Community Development

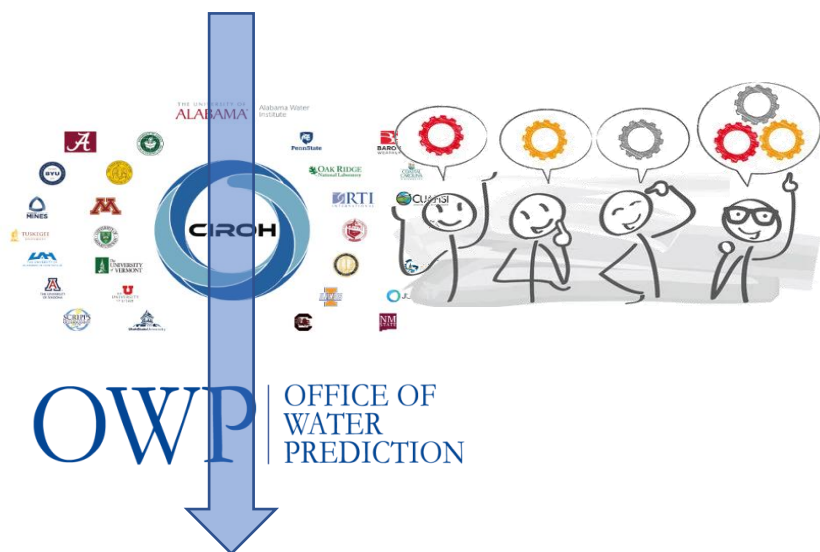


Need for Flexible, Scalable,
Exploratory Evaluation Tools



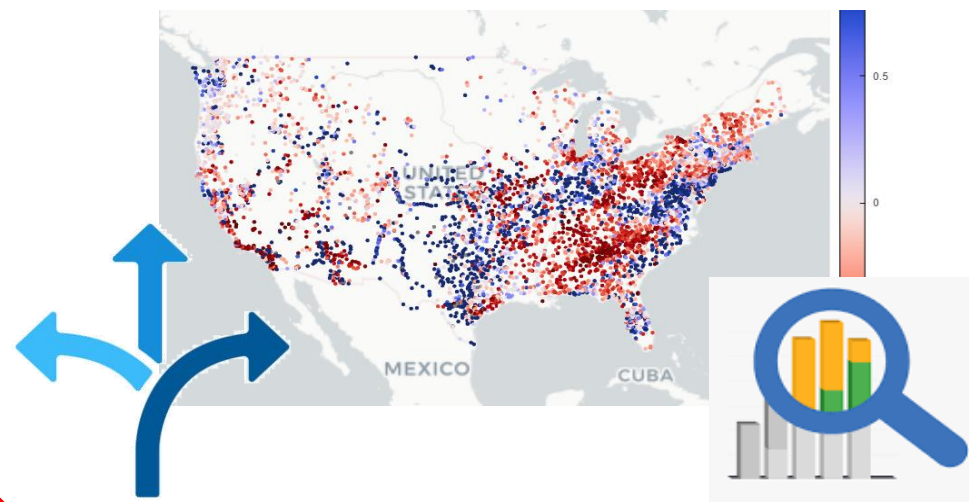
Two Parallel Goals

1) Establish **evaluation standards** in collaboration with OWP, CIROH and broader hydrologic forecasting community → How will we judge 'improvements'



Workshop Focus

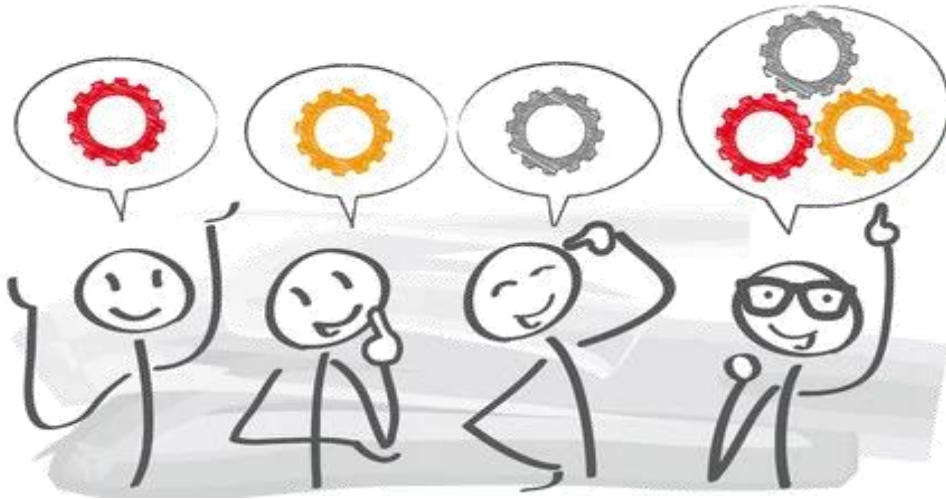
2) Create a **set of tools** for hydrologic model and forecast evaluation that are scalable and flexible for the wide range of use cases and users, that enable **highly exploratory evaluation**, and that foster open community development



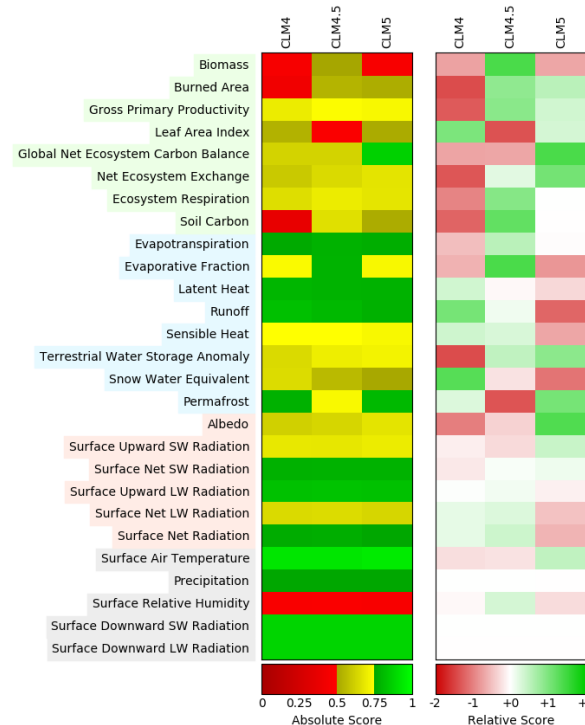
Envisioned Strategy for Goal #1

Collaboration across
CIROH and OWP

Converge on some standards for
performance assessment and
comparisons across research



e.g., CIROH Evaluation Working Group



ILAMB Example
<https://www.ilamb.org/>

Strategy for Goal #2



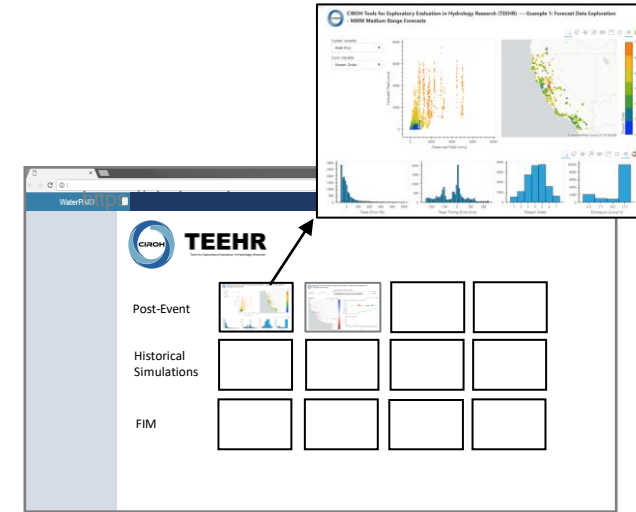
Tools for Exploratory Evaluation
in Hydrologic Research

- Still in early stages of development
- **Tiered** (see what we did there) levels of use, interaction and contribution
- Seeking feedback on needs for different CIROH projects and use cases (today and afterwards)

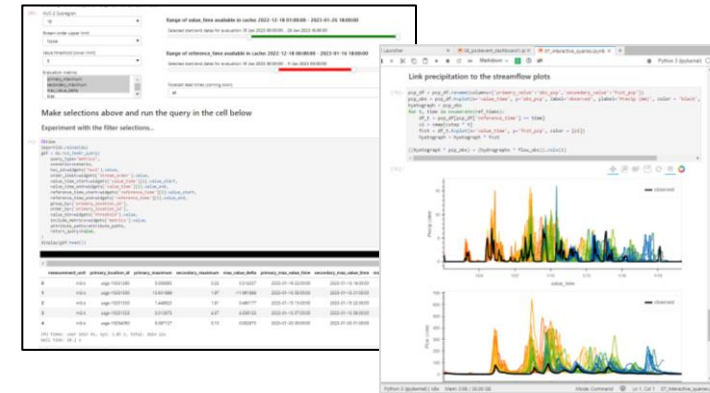
Tiered Levels of Use



Servable Dashboards



Notebook Templates

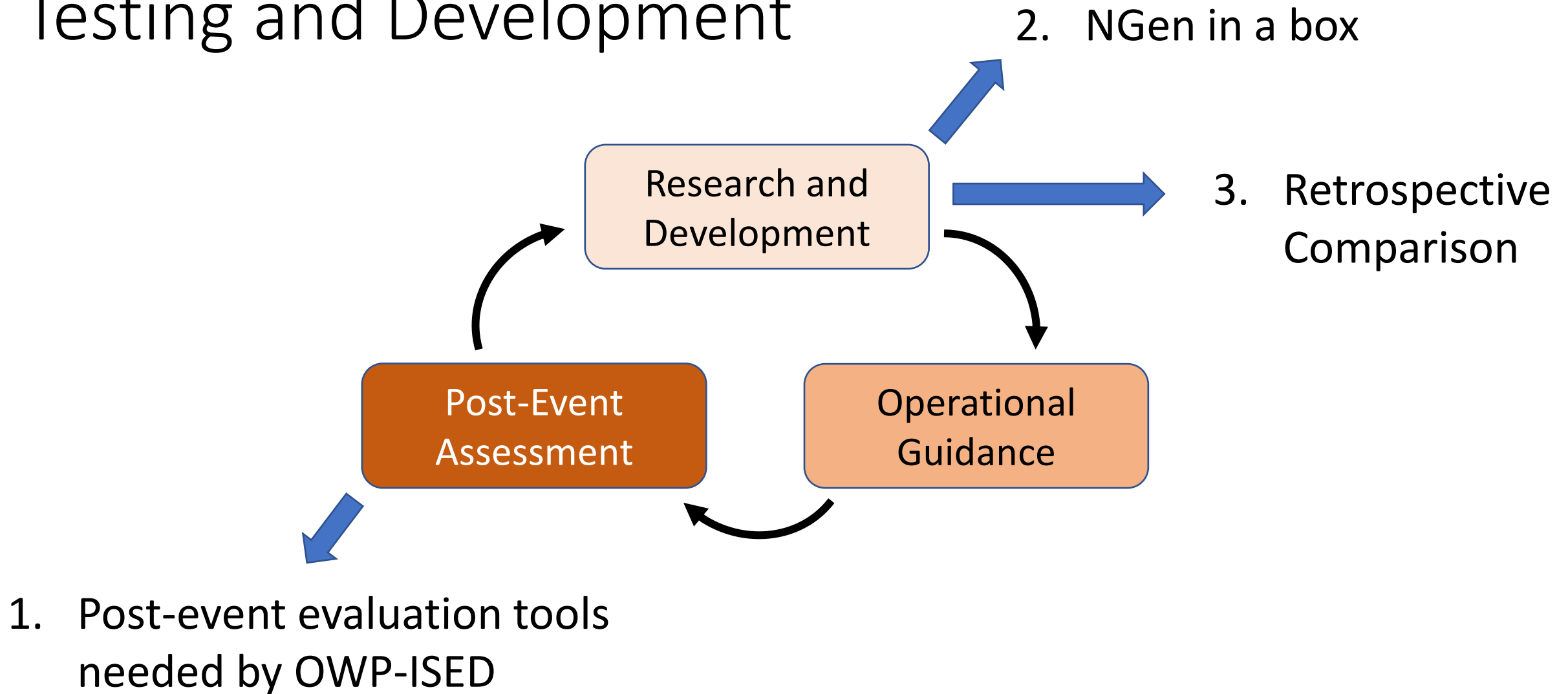


Python Package

```
qry = tod.get_metrics(  
    primary_filepath=USGS,  
    secondary_filepath=SHORT_RANGE,  
    crosswalk_filepath=CROSSWALK,  
    group_by=["primary_location_id", "reference_time"],  
    order_by=["primary_location_id"],  
    include_metrics="all",  
    filters=[  
        {  
            "column": "primary_location_id",  
            "operator": "=",  
            "value": "usgs-10336676"  
        },  
        {  
            "column": "reference_time",  
            "operator": "=",  
            "value": "2023-01-02 16:00:00"  
        }  
    ],  
    return_query=True  
)  
print(qry)
```

workshop

Initial Use Cases for Testing and Development

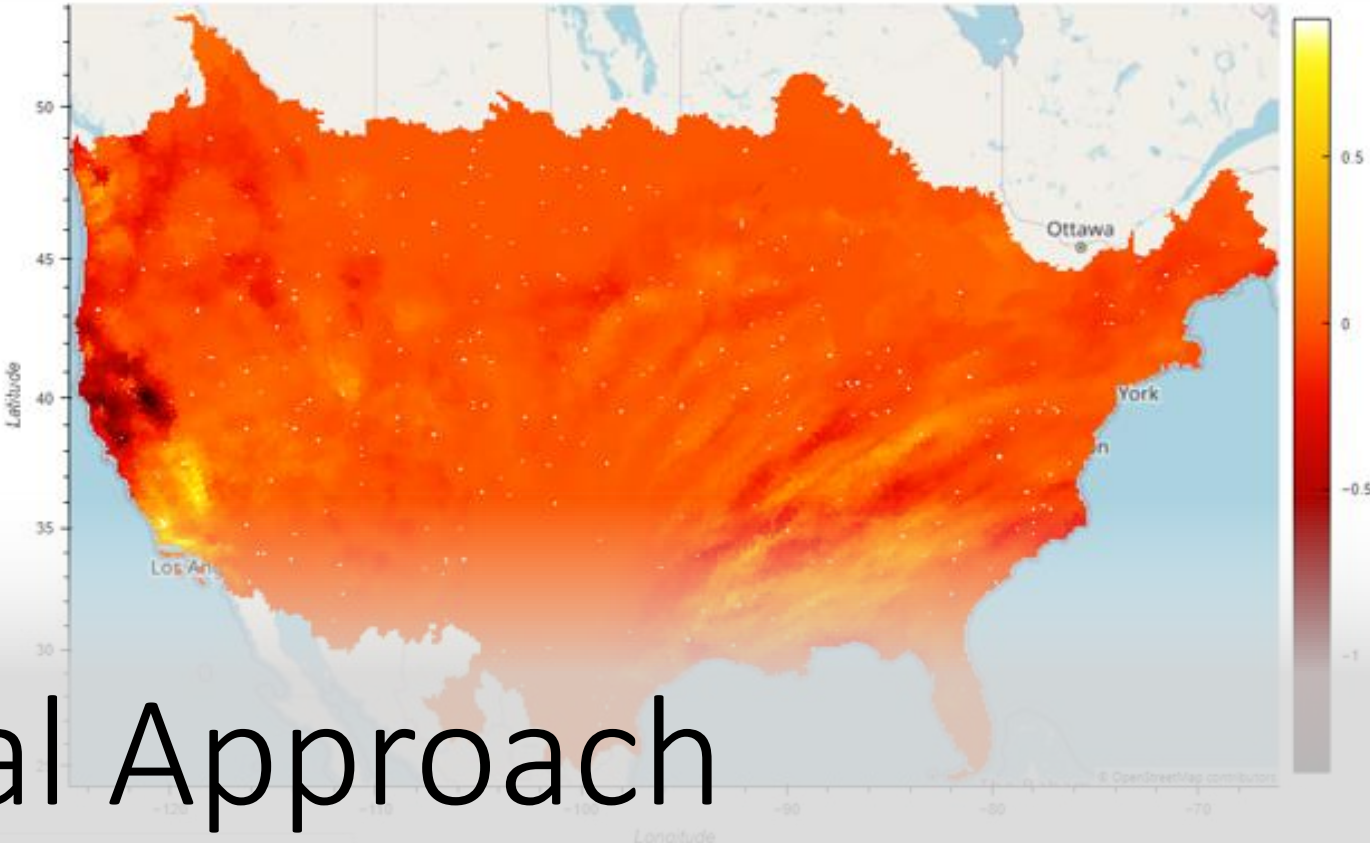


timeseries	
reference_time	[datetime]
value_time	[datetime]
value	[float]
measurement_unit	[string]
configuration	[string]
location_id	[string]

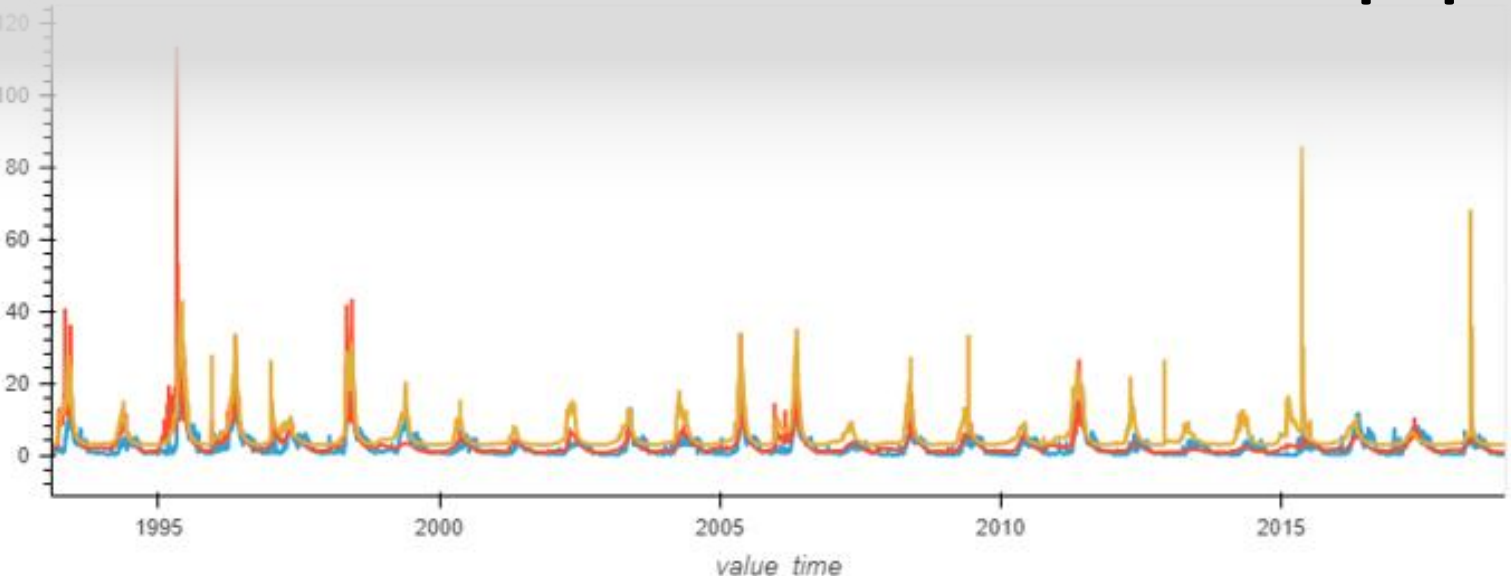
geometry	
id	[string]
name	[string]
geometry	[wkt or wkb]

attribute	
location_id	[string]
attribute_name	[string]
attribute_value	[string]
attribute_unit	[string]

crosswalk	
primary_location_id	[string]
secondary_location_id	[string]



Technical Approach



```

qry = tqd.get_metrics(
    primary_filepath=USGS,
    secondary_filepath=SHORT_RANGE,
    crosswalk_filepath=CROSSWALK,
    group_by=["primary_location_id", "reference_time"],
    order_by=["primary_location_id"],
    include_metrics="all",
    filters=[
        {
            "column": "primary_location_id",
            "operator": "=",
            "value": "usgs-10336676"
        },
        {
            "column": "reference_time",
            "operator": "=",
            "value": "2023-01-02 16:00:00"
        }
    ]
),
return query=True

```


Goals

Goal again: Create a **set of tools** for hydrologic model/forecast evaluation that are scalable and flexible for the wide range of use cases and users, that enable **highly exploratory evaluation**, and that foster open community development

Status: This is a CIROH funded research project. We are 9 months (~6 months of working time) into the project. We have a **prototype level** set of tools built and are figuring out what works and what doesn't. There is lots of work to do; we are open to feedback.

Objectives

- **Easy to use** tools that will form the backbone of exploratory forecast and simulation evaluation and visualization.
- Target audience is **data scientists, hydrologists, researchers**, maybe building a dashboard for a less technical audience.
- Use **familiar tooling** (Python, Pandas, Xarray, etc. -> Pangeo)
- Engage with the **community**; clean code, good documentation and well documented examples.
- Make the library **fast, efficient, scalable** to meet needs of different users and use cases.
- Not directly tied to any specific data source (NWM, USGS, etc.)
- There is so much great work by others; don't reinvent the wheel, keep a narrow focus and scope of our work

Technical Agenda Items

- Data Models
- Study Cache Structure
- Fetching and Loading Data
- Querying
- Visualization

Work in progress!

Loading

Too many to list...



netCDF



Zarr



dask



pandas

Cache/DB



Parquet

Query



DuckDB



pandas

Visualize



Panel



hvPlot



HoloViews



GeoViews



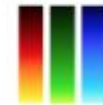
Datashader



Lumen



Param



Colorcet



Apache Parquet is an open source, column-oriented data **file format** designed for efficient data storage and retrieval. It provides efficient data compression and encoding schemes with enhanced performance to handle complex data in bulk.



DuckDB is an in-process SQL OLAP **database** management system.

With respect to TEEHR objectives, Parquet files and DuckDB provide a powerful way to store large amounts of timeseries data with a small storage footprint while still allowing relatively fast access via SQL queries. DuckDB can query Parquet files directly, including files stored in a cloud bucket.

Login to AWI CIROH JupyterLab



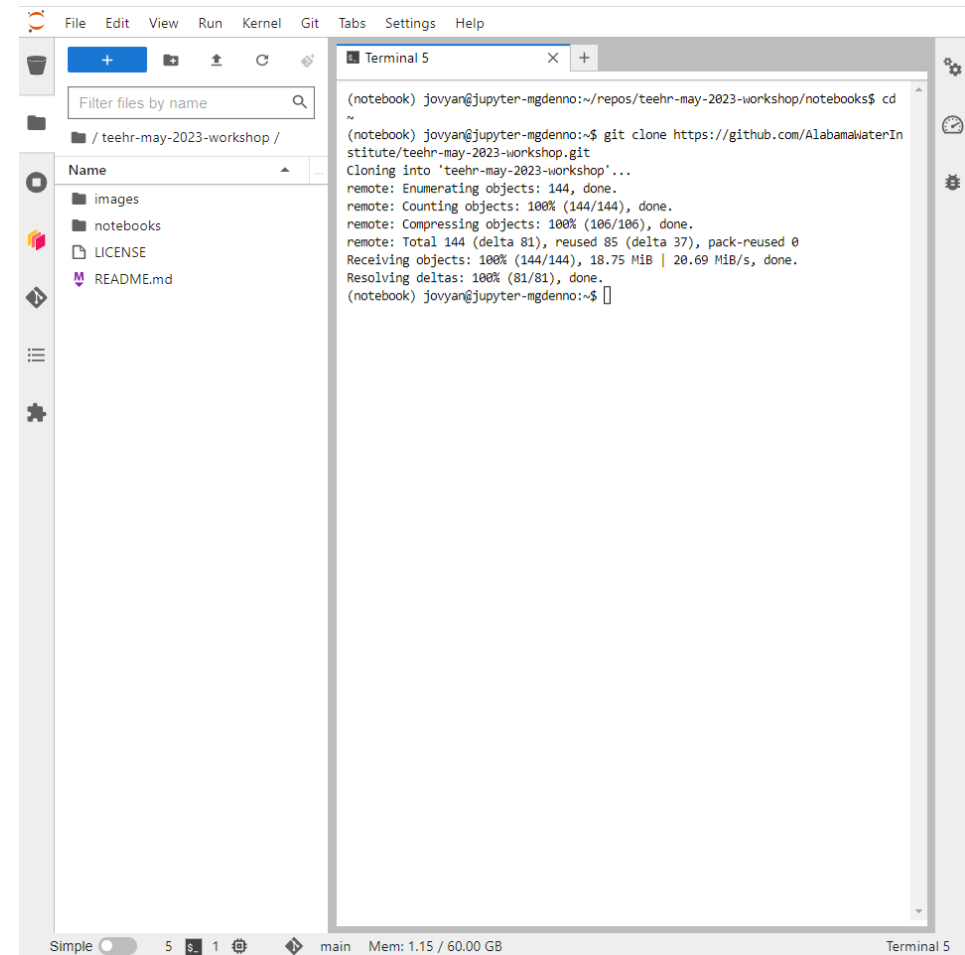
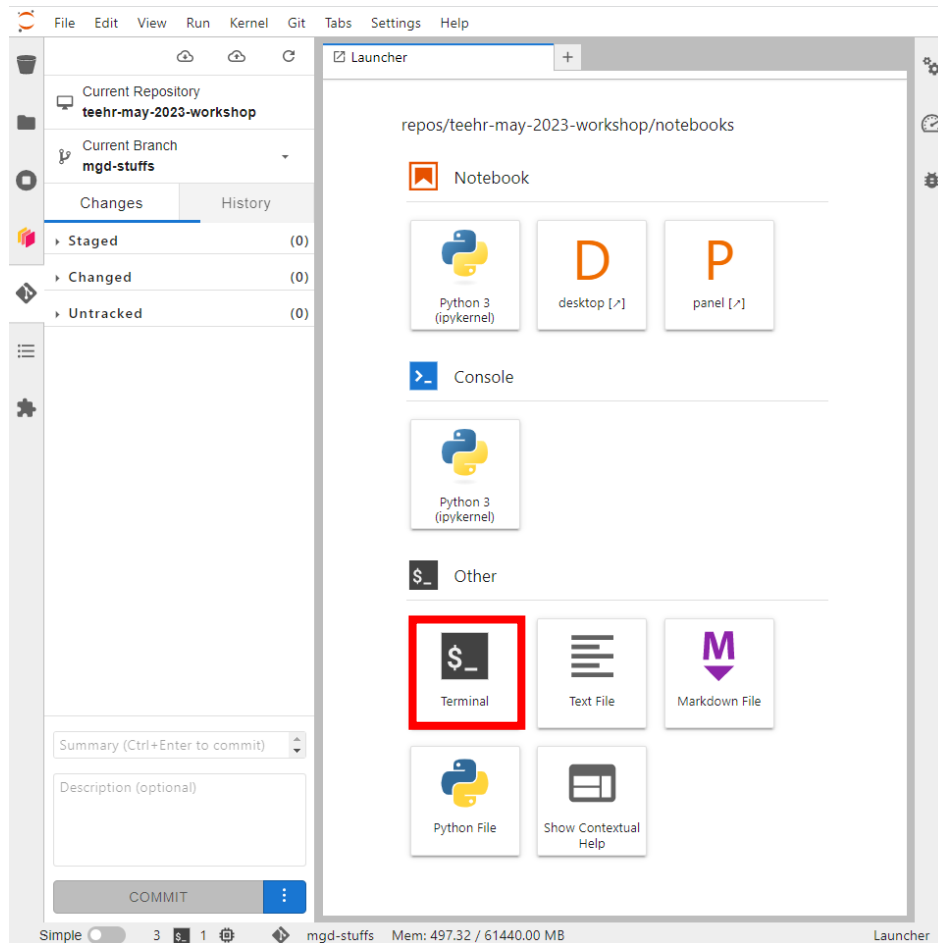
1. Login to your GitHub account
2. Accept invitation to join the **AlabamaWaterInstitute** GitHub organization (in GitHub)
3. Go to:
<https://ciroh.awi.2i2c.cloud/hub/login>
4. Authorize 2i2c-org to access your GitHub account.

1. Open “Launcher” if not already open:

i. File > New Launcher

2. Open “Terminal” and run:

 \$ git clone https://github.com/AlabamaWaterInstitute/teehr-may-2023-workshop.git



Future Work...too much to list

- Work with the community to identify standards, etc...
- Enhance the tools to work with cloud storage
- Add additional queries and metrics
- Access data remotely instead of downloading, query direct from bucket, etc.
- Research how to better utilize Parquet file strengths (i.e. wider tables)
- Tightly integrated visualization components
- Possibly prepare and stage commonly needed datasets, for example this could include, USGS gages and associated hydrologic attributes, HUC2-HUC12 polygons and weights files, etc.
- Tighter integration with NextGen
- Speed, speed, speed
- Data validation for the cache
- <https://github.com/RTIInternational/teehr/discussions/32>