# Meeting the Needs of Interdisciplinary Critical Zone Scientists by Leveraging and Linking Existing Domain Repositories

Jeffery Horsburgh[1], Kerstin Lehnert[2], Christopher Calloway[3], and Jerad Bales[4]

[1]Utah State University
[2]Columbia University
[3]University of North Carolina at Chapel Hill
[4]Consortium of Universities for the Advancement of Hydrologic Science, Inc.

November 21, 2022

## Abstract

Critical Zone (CZ) scientists study the system of coupled chemical, biological, physical, and geological processes operating together across all scales to support life at the Earth's surface (Brantley et al., 2007). In 2020, the U.S. National Science Foundation funded a new network of Thematic Cluster projects who are working collaboratively to answer scientific questions related to effects of urbanization on CZ processes; CZ function in semi-arid landscapes and the role of dust in sustaining these ecosystems; processes in deep bedrock and their relationship to CZ evolution; recovery of the CZ from disturbances such as fire and flooding; and changes in the coastal CZ related to rising sea level. Given the diversity of data being collected by these projects, supporting data collection, access, and archival for the larger network presents significant challenges. Leveraging existing repositories and cyberinfrastructure provides many benefits, but still poses the questions of which repositories to use and how to enable discovery of and access to data that may be deposited across different repositories. This presentation describes new cyberinfrastructure development that leverages existing, domain-specific data repositories to enable managing, curating, disseminating, and preserving data from the new network of CZ Thematic Cluster projects. A distributed architecture is under development that links existing data facilities and services, including HydroShare, EarthChem, SESAR, and eventually other systems as needed, via a CZ Hub that provides tools for simplified data submission, discovery and access, and links to computational resources for data analysis and visualization in support of CZ synthesis efforts. Our goal is to make data, samples, and software collected by the Thematic Cluster projects Findable, Accessible, Interoperable, and Reusable (FAIR), using existing domain-specific repositories. This collaboration among repositories to deliver integrated data services for an interdisciplinary science program may provide a template for future development of integrated, interdisciplinary data services. Brantley, S.L., M.B. Goldhaber, V. Ragnarsdottir (2007). Crossing disciplines and scales to understand the Critical Zone. Elements 3, 307-314, doi:10.2113/gselements.3.5.307.

# Meeting the Needs of Interdisciplinary Critical Zone Scientists by Leveraging and Linking Existing Domain Repositories

**Jeffery S. Horsburgh**

Utah State University

**Kerstin Lehnert, Chris Calloway, Jerad Bales**

# Critical Zone Collaborative Network

- In 2020 NSF funded the next phase of their Critical Zone research program

- Nine Thematic Cluster study areas with a wide range of geological, climatic, and land use settings working to better understand the evolution and function of the Critical Zone

- One Coordinating Hub to help coordinate activities across Clusters – including data management

**BEDROCK**
Expanding knowledge of the deep critical zone and its feedbacks with surface processes.

**COASTAL**
Investigating the processes that transform landscapes and fluxes between land and sea.

**DYNAMIC WATER**
Advancing the understanding of the interactions among dynamic water storage, CZ processes, and water provisioning in western U.S. montane ecosystems.

**BIG DATA**
Using field observations, existing data, & advanced statistical and process-based tools to investigate how the Critical Zone responds to disturbances.

**DRYLANDS**
Quantifying and predicting dryland carbon budgets across land-use and climatic gradients.

**GEOMICROBIO**
Studying how soil microbes, roots, mineral composition, and soil organic matter interact and drive Critical Zone biogeochemistry and soil formation.

**CINET**
Investigating the role of critical interfaces for regulating the storage & transport of material such as water, sediment, carbon, & nutrients.

**DUST^2**
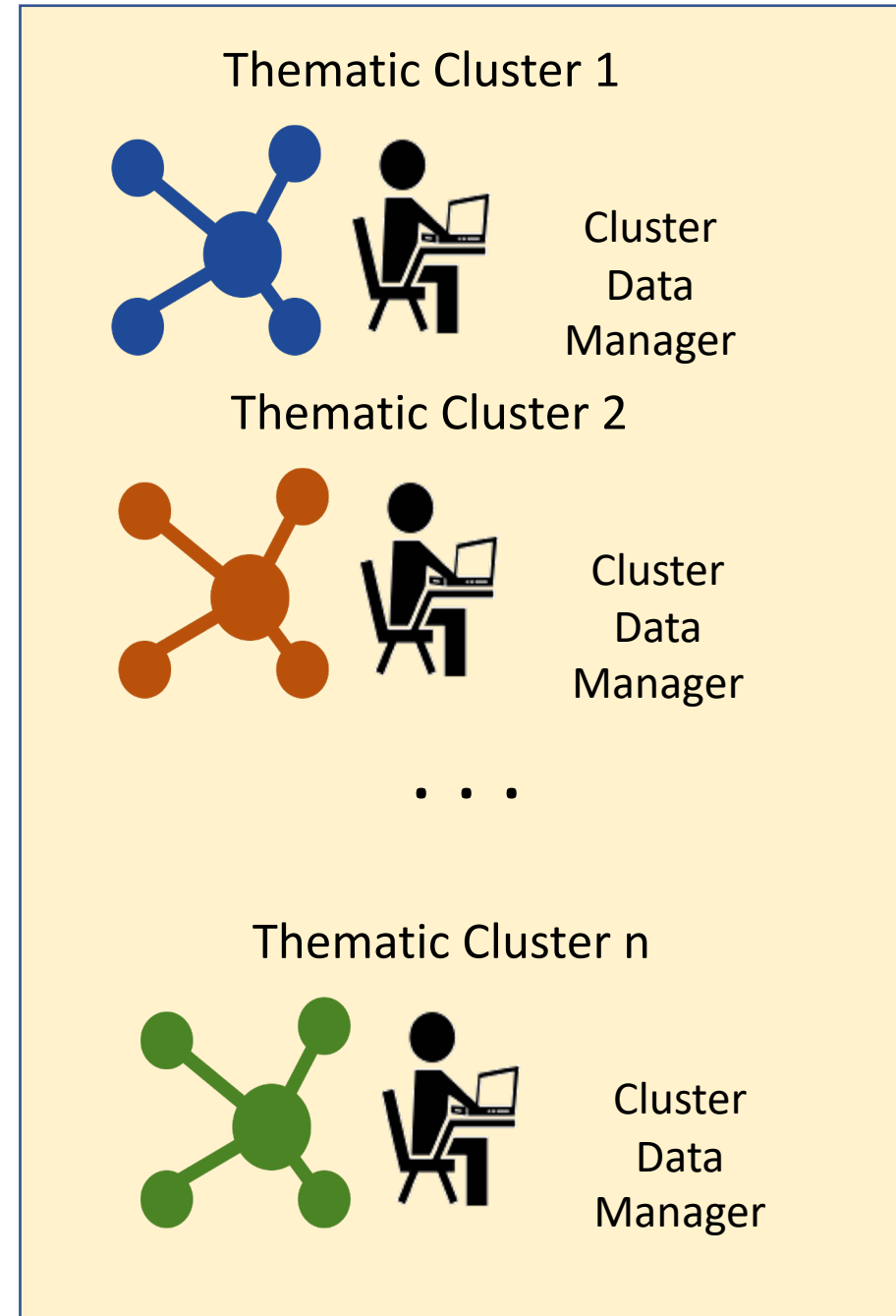A source-to-sink investigation of the dust system in the southwestern US as a component of the critical zone.

**URBAN**
Studying the interaction between the geologic template and the urban footprint and the effects on critical zone processes along the Eastern Seaboard.

# Thematic Cluster Projects

- Activities
  - Data collection
  - Data aggregation
  - Data QA/QC
  - Defining data products
  - Metadata creation
- Data Management Plans
  - Each Thematic Cluster submitted their own Data Management Plan
  - Flexibility for local data management

Thematic Cluster 1

Cluster Data Manager

Thematic Cluster 2

Cluster Data Manager

. . .

Thematic Cluster n

Cluster Data Manager

# Challenges

- Thematic Cluster teams and data are diverse

- Some are collecting new data, others are aggregating existing data, some are doing both

- No single data repository will meet the needs of interdisciplinary Critical Zone Scientists

CZ Hub Objective: Provide a robust cyberinfrastructure for **F**indable, **A**ccessible, **I**nteroperable, and **R**eusable (FAIR) data from the CZ Net Thematic Clusters

| Thematic Clusters | |
|---|---|
| 1 | Bedrock |
| 2 | Coastal |
| 3 | Dynamic Water |
| 4 | Big Data |
| 5 | Drylands |
| 6 | Geomicrobio |
| 7 | CINet |
| 8 | Dust^2 |
| 9 | Urban |

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.

# CZ Hub Approach

- Link existing data facilities and services, including:
  - HydroShare
  - EarthChem
  - System for Earth Sample Registration (SESAR)
  - OpenTopography
  - Other repositories, as needed
- Develop a central CZ Hub that provides
  - Services for easy data submission
  - Integrated data discovery and access

# CZ Hub Approach

Diverse data and research products from CZ scientists

New Submission Portal promoting repository selection, metadata, templates, and formats

Using existing Earth science data repositories via automated submission

Making products Findable, Accessible, Interoperable, and Reusable (FAIR)



**CZ Thematic Cluster Network**
- Data collection
- Data aggregation
- Local data management
- Quality assurance/quality control
- Metadata creation

Cluster Data Manager

Cluster Data Manager

Cluster Data Manager

**Data Submission Portal**
- User support specialist
- Data curation support
- Documentation

**Data Submission**
- **Metadata Templates**
- **Data Format Standards**
- **Controlled Vocabularies**
- **Data Upload Templates**
- **Sample Registration**
- **Unique Identifier Management**

**Computational Resources**
- Jupyter Notebook support
- Computation and modeling
- Synthesis studies
- Reproducible analyses

**Repositories for Data and Research Products**
- Permanent data archival and publication
- Access control for embargoed data
- Open access for public datasets
- Citable data

API EarthChem
- Samples
- Geochemistry
SESAR SYSTEM FOR EARTH SAMPLE REGISTRATION

API HYDROSHARE
- Time series
- Geospatial
- NetCDF
- Models
- Generic

API OpenTopography
- LiDAR
- High resolution topography
- Elevation products

API **Other Repositories**

**Catalog Services**
- Cross-repository view of CZ data and research products
- Discovery based on authors, keywords, geograhic area, time, Cluster
- Schema.org metadata implementation

# Data Submission Portal

- New, web application to support CZ Net

- Enables submission to multiple geoscience data repositories through one portal

- Getting data to the right repository

- Submission directly through portal

Empower data managers to curate research products within appropriate repositories with support from our team

# User Account Management

- Authentication and user management via ORCID

- Using an account that most people have already

- Enables authentication across multiple repositories

# Supported Repositories

- Operate and partner with existing repositories
  - Promote the use of FAIR principles
  - Permanent data archival and publication
  - Access control for embargoed data
  - Open access for public datasets
  - Citable data
  - Leverage existing NSF investment in CI

# Which repository?

- Which repositories to target for different data types?

- Assist data managers in selecting an appropriate repository
  - Geospatial data
  - Data derived from physical samples
  - Hydrologic time series
  - Data and code packages
  - Models

# Step1: Log in

- User logs into the Portal using their ORCID
- User accounts in the Portal are associated with the ORCID

# Step 2: Choose Repo

- User chooses a repository to submit to from the Submit Data page
- The user authorizes the Portal to submit to HydroShare
- The one-time authorization is stored in the user's profile

# Step 3: Create Content

- User enters metadata and selects content files on the data submission form for the chosen repository
- Each repository has its own submission form
- Submission forms are built from a JSON schema that defines:
  - Required and optional metadata
  - Default values
  - Etc.

# Step 4: Submit Content

- Metadata and data files are sent to the repository
- A new resource is created in the repository
- A new record is created in the user's My Submissions page
  - Return later to Edit
  - Export submissions to file

# JSON Schema-based Metadata

- A JSON schema defines required and optional metadata for each repository
- Submissions validated based on JSON schema
  - Data types
  - Default values
  - Required/optional
- Data submission form dynamically built from the JSON schema
- Adding a new repository to the Portal means adding a new JSON schema

```
{
  "title": "Resource Metadata",
  "description": "A class used to represent the metadata for a resource",
  "type": "object",
  "properties": {
    "title": {
      "title": "Title",
      "description": "A string containing the name given to a resource",
      "maxLength": 300,
      "type": "string"
    },
    "abstract": {
      "title": "Abstract",
      "description": "A string containing a summary of a resource",
      "type": "string"
    },
    "language": {
      "title": "Language",
      "description": "A 3-character string for the language in which the metadata and content of a resource
      "type": "string"
    },
    "subjects": {
      "title": "Subject keywords",
      "description": "A list of keyword strings expressing the topic of a resource",
      "default": [],
      "type": "array",
      "items": {
        "type": "string"
      }
    },
    "creators": {
      "title": "Creators",
      "description": "A list of Creator objects indicating the entities responsible for creating a resource"
      "default": [],
      "type": "array",
      "items": {
        "$ref": "#/definitions/Creator"
      }
    },
    "contributors": {
      "title": "Contributors",
      "description": "A list of Contributor objects indicating the entities that contributed to a resource",
      "default": [],
```

# JSON Schema-based Metadata

- Data submission form dynamically built from the JSON schema
- Files and metadata sent directly to repository
- Data Submission Portal maintains a record of submission

# Promoting best practices

- Repository functionality is not specific to a community of users
- CZ Net may want to use:
  - Community standard formats
  - Templates
  - Best practices
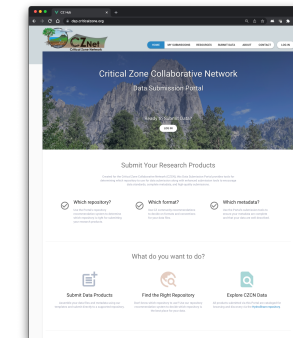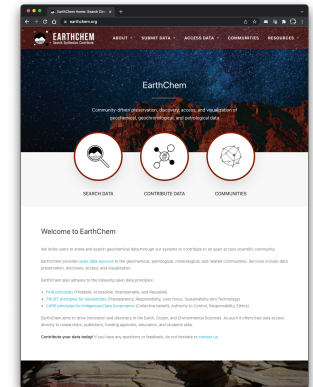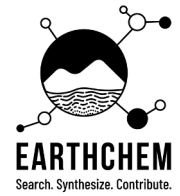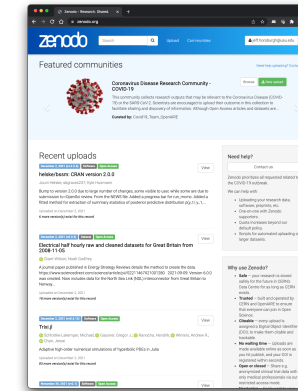- We can promote those through the Data Submission Portal

# Without the DSP

- Data managers must navigate user interfaces of multiple systems

- Must keep track of what has been submitted to each one

- Difficult for CZ Hub Team to track what has been submitted

# With the DSP

- Data managers need only interact through one user interface

- Submissions to all repositories tracked in one place

- Submissions automatically registered for cataloging/discovery

# Data Submission Portal Advantages

- Through validation, promote consistency in CZ Net data products across repositories
- Ensure data products end up in an appropriate, trusted repository
- Enforce minimum metadata requirements
  - Consistent keywords
  - Funding agency/grant information
- Enable use of controlled vocabularies where needed
- Promote templates, common formats, and best practices
- Enable data managers to use a single interface/tool to submit data
- Enable simple and consistent registration of CZ Net datasets with a metadata index for discovery
- Helping Thematic Clusters keep track of what has been submitted

# CZ Net Catalog Services

- Cross-repository view of CZ Net data and research products

- Discovery based on authors, geographic area, time, cluster

- Schema.org metadata

A coordinated view and data discovery service(s) for all the data produced within the collaborative network to ensure that data are **Findable** and **Accessible**.

HydroShare datasets discoverable via Google Dataset Search

# CZ Net Catalog Services

CZO "Community" in HydroShare with individual "Groups" for each observatory

- Cross-repository view of CZ Net data and research products

- Discovery based on authors, geographic area, time, cluster

- Schema.org metadata

- Communities and Groups in HydroShare

A coordinated view and data discovery service(s) for all the data produced within the collaborative network to ensure that data are **Findable** and **Accessible**.