

# Enabling Scientific Discovery: Harnessing the Power of the National Science Data Fabric for Large-Scale Data Analysis (Session I & II)

**Presenter:** Michela Taufer<sup>1</sup>

**Other contributors:** Heberth Martinez<sup>1</sup>, Paula Olaya<sup>1</sup>, Asshish Panta<sup>2</sup>, Amy Gooch<sup>2</sup>, Jack Marquez<sup>1</sup>, Gabriel Laboy<sup>1</sup>, Jay Ashworth<sup>1</sup>, Giorgio Scorzelli<sup>2</sup> and Valerio Pascucci<sup>2</sup>

<sup>1</sup>University of Tennessee Knoxville, <sup>2</sup>University of Utah



# Acknowledgments

The authors of this tutorial would like to express their gratitude to:

- NSF through awards 2138811, 2103845, 2334945, 2138296, and 2331152
- [Dataverse](#)
- [Seal Storage](#)
- [Rodrigo Vargas](#), Vargas Lab, University of Delaware
- Werner Sun, [CHESS](#), Cornell University
- DOE SBIR Phase II award DE-SC0017152

Any opinions, findings, conclusions, or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.



# Schedule

The half-day tutorial is organized into four sessions:

## Session I (30 mins):

This session begins with an overview of the NSDF and addresses users' challenges identified through interviews.

## Session II (1 hour):

This session offers a hands-on experience with NSDF services, focusing on visualization and dashboard creation for Earth science datasets.

## Session III (1 hour):

This session delves deeper into NSDF services tailored for the management and analysis of datasets exceeding 100TB.

## Session IV (30 mins):

This session concludes with an interactive Q&A, allowing attendees to discuss applications of NSDF in various research fields.

# Prerequisites



Tutorial

## Step 0: Access to GitHub

To run this tutorial, you need to have a GitHub account.

- You can create one following the instructions here:

<https://docs.github.com/en/get-started/start-your-journey/creating-an-account-on-github#>

- Now you can login into GitHub

<https://github.com/login>

## Step 1: Create Codespaces

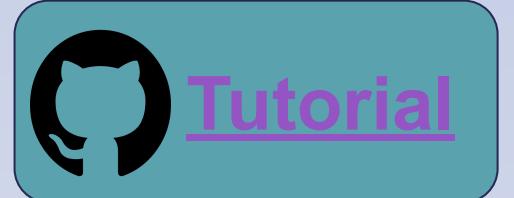
Use your GitHub account to run this tutorial with GitHub codespaces

- Access this link:  
[NSDF Tutorial 2024](#)
- Click on green button  
“Create codespace”

Create codespace

✓ Image found.  
⠼ Building container...

# Creating the GitHub Codespace



Create a new codespace

**Repository**  
To be cloned into your codespace  
nsdf-fabric/Tutori... ▾  
Codespace usage for this repository is paid for by jackdmarquez.

**Branch**  
This branch will be checked out on creation  
main ▾

**Dev container configuration**  
Your codespace will use this configuration  
NSDF Tutorial - Session II ▾

**Region**  
Your codespace will run in the selected region  
US East ▾

**Machine type**  
Resources for your codespace  
2-core ▾

Click this button -> **Create codespace**

→ Let's start:

[NSDF Tutorial 2024](#)



# Tutorial Goals

This tutorial demonstrates end-to-end analysis of scientific data through NSDF services

## Tutorial Goals

**Construct a modular workflow** that combines your application components with NSDF services

**Upload, download, and stream data** to and from **public and private storage** solutions

**Deploy the NSDF dashboards** for large-scale **data access, visualization, and analysis**



National Science Data Fabric



[www.sci.utah.edu](http://www.sci.utah.edu)



THE UNIVERSITY OF  
TENNESSEE  
KNOXVILLE



Powered by ViSUS



SDSC

IBM





# Session I: Understanding and Addressing User's Pain Points

Surveying Community Needs and Realities



National Science Data Fabric



[www.sci.utah.edu](http://www.sci.utah.edu)



THE UNIVERSITY OF  
TENNESSEE  
KNOXVILLE



Powered by ViSUS



SDSC

IBM



# Identifying Pain Points: User Interviews



## Identify Users

- Diverse roles: Domain scientists, CI professionals, developers
- Diverse domains: materials science, climate, earth sciences, astronomy, and more!
- Diverse institutions: R1 universities, teaching colleges, MSIs, national labs, experimental facilities

## Target Questions

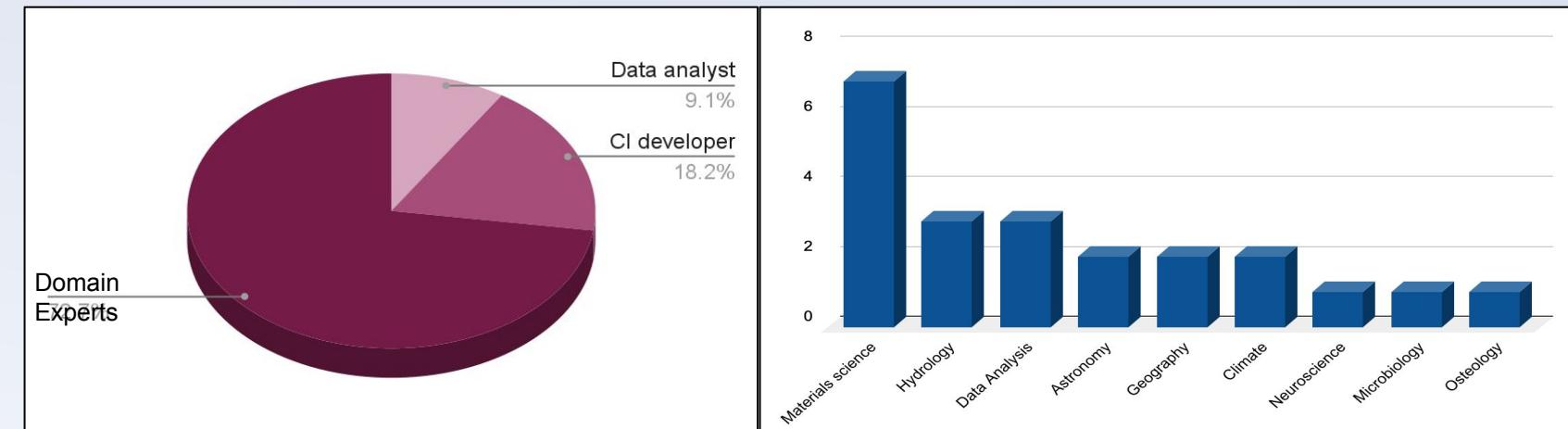
- General questions about data storage, data form, metadata storage, WMS, data catalogs, programming languages
- Specific questions about unique challenges related to role, domain, and institution

## Analyze Results

- Identify cross-cutting concerns
- Identify concerns consistent for roles, domains, and institutions

## Diagnose Pain Points

- Distill concerns into concrete problem statements
- Translate into objectives, actionable items, and milestones



# Identifying Pain Points: Testimonials



*"If we make it easy, people will share their data. We need to extend the scalability of our infrastructure through community supplied storage." - Project PI, Materials Science*

*"We don't have a plan to scale storage if our cyberinfrastructure takes off." - Cyberinfrastructure developer, Hydrology*

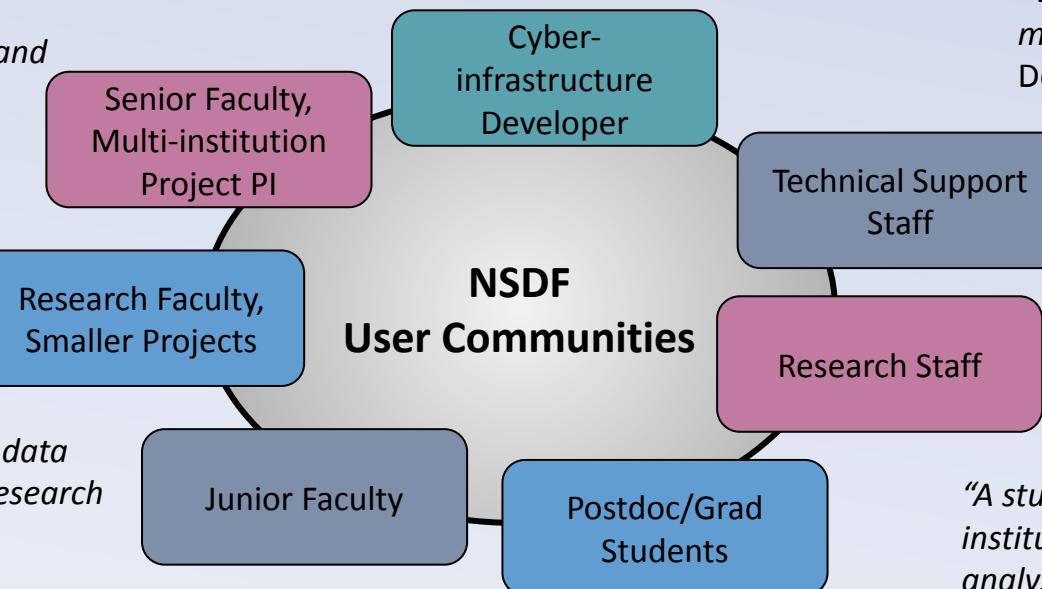
*"We can't scale to PB data without massive funding or infrastructure. We need centralized point of access to federated data" - Cyberinfrastructure Developer, Materials Science*

*"The time and effort for using public repositories, and limited realized gains limits our data sharing" - Faculty, Materials Science*

*"Remote quality control during acquisition would let us better use beamline time" - Faculty, Materials Science*

*"Before I had funding to run my own experiments, data shared by a friend at a national lab launched my research career" - Faculty, Materials Science*

*"I'm perplexed by the lack of urgency around reproducible and replicable processes for data management" - Faculty, Geography*



*"Our old data is kept on (external) drives. It's hard to keep things organized" - Graduate Student, Materials Science*

*"Our long-term storage is a shelf of external hard drives" - Research Staff, Neuroscience*

*"We move data (from light source) by flying back with TB hard drives" - Graduate Student, Materials science*

*"We lack personnel to do basic development and maintenance of our systems" - Cyberinfrastructure Developer, Astronomy*

*"We can't hire enough system maintainers and have research funding" - Data Analyst Group Lead*

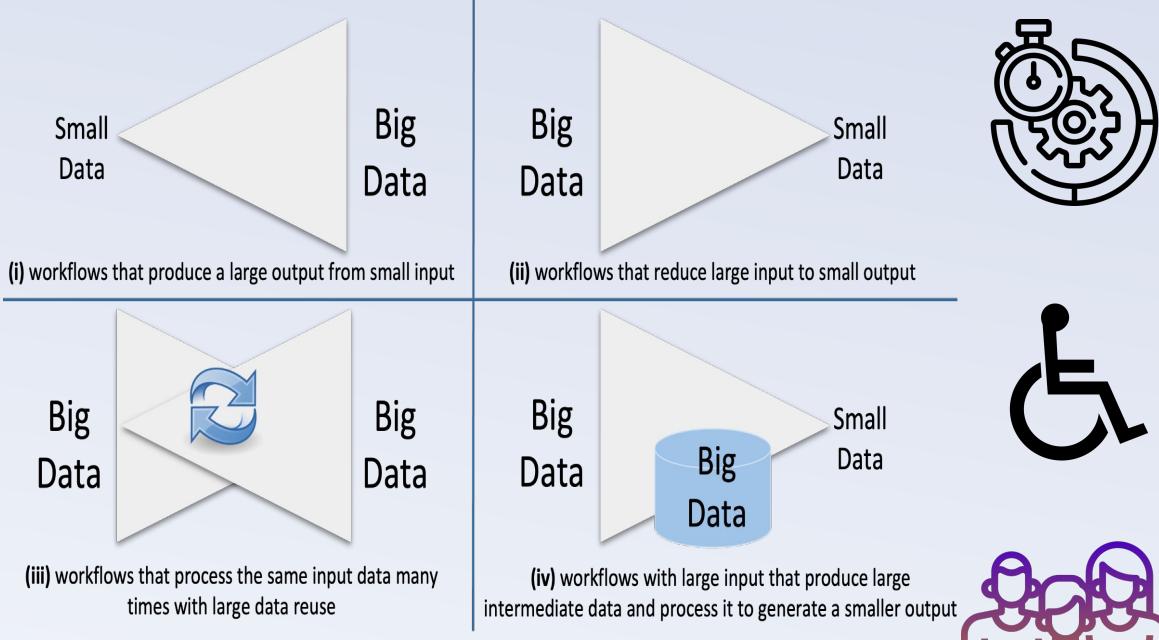
*"Jupyter notebook access to data (TBs) would reduce barrier to entry" - Research Staff, Climate Science*

*"A student copies GBs of data from the scientist to my institution. I download to my laptop to prototype analysis. It is cumbersome and limits testing." - Research Staff, Data Analysis*



# Pain Points Inform Data Science Strategies

NSDF focuses on the main classes of workflows that **CHALLENGE** data-intensive scientific investigations



**Scalability** Software stack scale from leadership computing to commodity hardware (even handheld)

**Information Streams** federated resources that minimize data transfers with trustworthy sharing of information

**Resource Efficiency** Allows teams to work effectively with limited access to human and physical resources

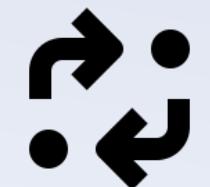
**Data Management** Standardized data and metadata management tools avoid replicated work

**Accessibility** Facilitate data-sharing processes for open and secure environments

**Timeliness** Immediate access and use of remote information without bulk data transfers

**Workforce Development** Trained of CI professionals

**Replicability** Programs/data versioning with FAIR identifiers throughout the scientific investigation



10



# What do Scientists Need?

Information streaming:  
knowledge exploration  
at run time



Access and use of data from  
remote facilities  
(experimental, computing,  
and repository)

Findable, Accessible,  
Interoperable, Reusable  
(FAIR) data

Equity in  
access and use

Remote queries with minimal  
data movement to mitigate  
bandwidth and storage  
issues



# Session I: Implementing an Accessible & Tightly Integrate Data Fabric

Designing, developing, and deploying equitable services



We are building a **holistic ecosystem to democratize data-driven scientific discovery** by connecting an open network of institutions, including minority-serving institutions, with a shared, modular, containerized data delivery environment.





100G Core

Terabit Core



NSDF EntryPoints



OSG StashCaches



Both

Institutions and  
universities  
with resources  
to share



Today  
Aug 6 2022  
19:51:04 UTC

20:00:00 UTC

Aug 9 2022 00:00:00 UTC

Aug 9 2022 04:00:00 UTC

Aug 9 2022 08:00:00 UTC

Aug 9 2022 12:00:00 UTC

Aug 9 2022 16:00:00 UTC

Full screen





100G Core

Terabit Core



NSDF EntryPoints



OSG StashCaches



Both

Initiative to  
integrate  
minority  
serving  
institutions



20:00:00 UTC

Aug 9 2022 00:00:00 UTC

Aug 9 2022 04:00:00 UTC

Aug 9 2022 08:00:00 UTC

Aug 9 2022 12:00:00 UTC

Aug 9 2022 16:00:00 UTC

Full screen

CESIUM ion Upgrade for commercial use. Data attribution



100G Core

Terabit Core



NSDF EntryPoints



OSG StashCaches



Both

Initiative to  
integrate  
large scale  
scientific  
projects



20:00:00 UTC

Aug 9 2022 00:00:00 UTC

Aug 9 2022 04:00:00 UTC

Aug 9 2022 08:00:00 UTC

Aug 9 2022 12:00:00 UTC

Aug 9 2022 16:00:00 UTC

Full screen

CESIUM ion Upgrade for commercial use. Data attribution



100G Core

Terabit Core



NSDF EntryPoints



OSG StashCaches



Both

# Initiative to integrate HPC resources



20:00:00 UTC

Aug 9 2022 00:00:00 UTC

Aug 9 2022 04:00:00 UTC

Aug 9 2022 08:00:00 UTC

Aug 9 2022 12:00:00 UTC

Aug 9 2022 16:00:00 UTC

Full screen

CESIUM ion Upgrade for commercial use. Data attribution



100G Core

Terabit Core



NSDF EntryPoints



OSG StashCaches



Both

# Initiative to integrate public cloud resources



Today

Aug 9 2022

19:51:04 UTC

20:00:00 UTC

Aug 9 2022 00:00:00 UTC

Aug 9 2022 04:00:00 UTC

Aug 9 2022 08:00:00 UTC

Aug 9 2022 12:00:00 UTC

Aug 9 2022 16:00:00 UTC

Full screen

Aug 9 2022

19:51:04 UTC

20:00:00 UTC

Aug 9 2022 00:00:00 UTC

Aug 9 2022 04:00:00 UTC

Aug 9 2022 08:00:00 UTC

Aug 9 2022 12:00:00 UTC

Aug 9 2022 16:00:00 UTC

Aug 9 2022

19:51:04 UTC

Aug 9 2022 00:00:00 UTC

Aug 9 2022 04:00:00 UTC

Aug 9 2022 08:00:00 UTC

Aug 9 2022 12:00:00 UTC

Aug 9 2022 16:00:00 UTC

Aug 9 2022

19:51:04 UTC

Aug 9 2022 00:00:00 UTC

Aug 9 2022 04:00:00 UTC

Aug 9 2022 08:00:00 UTC

Aug 9 2022 12:00:00 UTC

Aug 9 2022 16:00:00 UTC

Aug 9 2022



100G Core

Terabit Core



NSDF EntryPoints



OSG StashCaches



Both

# Initiative to integrate enterprise cloud and storage resources



Today  
Aug 9 2022  
19:51:04 UTC

20:00:00 UTC

Aug 9 2022 00:00:00 UTC

Aug 9 2022 04:00:00 UTC

Aug 9 2022 08:00:00 UTC

Aug 9 2022 12:00:00 UTC

Aug 9 2022 16:00:00 UTC

Full screen



100G Core  
Terabit Core

NSDF EntryPoints

OSG StashCaches

Both

A data fabric must be accessible and tightly integrated to coordinate data movement and use among geographically distributed teams or organizations



MINORITY SERVING-CYBERINFRASTRUCTURE CONSORTIUM





# National Science Data Fabric



National Science Data Fabric



# National Science Data Democratization Consortium: Engaging Industry Partners



MINIO



SEAL



WEKA



CLOUDFLARE



IBM Cloud



DoubleCloud



ALLUXIO



National Science Data Fabric



www.sci.utah.edu



THE UNIVERSITY OF  
TENNESSEE  
KNOXVILLE



ViSiCAR  
Powered by VISUS



JOHNS HOPKINS  
UNIVERSITY

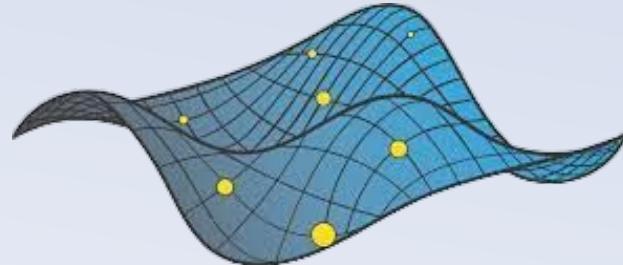
SDSC

IBM



22

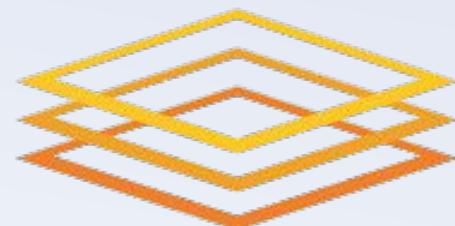
# Partnering with Existing NSF Initiatives



FABRIC



open storage network



Open Science Grid



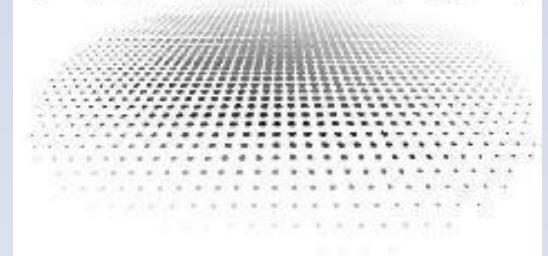
SDSC

IBM



23

NATIONAL DATA  
PLATFORM





# Partnerships with DoE Labs

- **Sandia National Laboratories**
  - Workflow containerization  
(Trustworthy Computing; Data Democratization)
- **Lawrence Livermore National Laboratory**
  - Thicket project (Large Scale Computing and Performance)
  - Flux project (Scheduling and Resource Management)
  - Fractale (Convergence of HPC, Cloud, and Edge)



National Science Data Fabric



[www.sci.utah.edu](http://www.sci.utah.edu)



THE UNIVERSITY OF  
TENNESSEE  
KNOXVILLE



Powered by VISUS



JOHNS HOPKINS  
UNIVERSITY

SDSC

IBM.





# Session I: Our Services and Successful Stories

Democratizing Access and Use of Large-scale Data

1



National Science Data Fabric



SDSC

IBM



25



100G Core

Terabit Core

NSDF EntryPoints

OSG StashCaches

Both

A **data fabric** must be **accessible** and **tightly integrated** to coordinate data movement and use among geographically distributed teams or organizations

Develop a **FAIR, AI-ready, transdisciplinary software stack** that is **easy to use, integrate, and scale**

Develop a federated data fabric: a suite of equitable **network, computing, and storage services** interoperating across the academic and commercial cloud  
(AWS, Azure, ...)



# NSDF Services

NSDF  
PLATFORM

Networking & Transfer



- NSDF-Plugin
- NSDF-Data Transfer

Computing & Monitoring



- NSDF-Cloud
- NSDF-Monitoring
- NSDF-Workflow

- NSDF-Fuse
- NSDF-Catalog
- NSDF-Stream
- NSDF-Dashboards
- NSDF-OpenVisus
- NSDF-Notebooks



Data & Storage



# NSDF Services

NSDF  
PLATFORM

Networking & Transfer



- NSDF-Plugin
- NSDF-Data Transfer

Computing & Monitoring



- NSDF-Cloud
- NSDF-Monitoring
- NSDF-Workflow



- NSDF-Fuse
- NSDF-Catalog
- NSDF-Stream
- NSDF-Dashboards
- NSDF-OpenVisus
- NSDF-Notebooks

Data & Storage



# Section I: Sharing Use-Inspired Research Stories

Decentralizing Research Hubs for Transformative Scientific Discovery

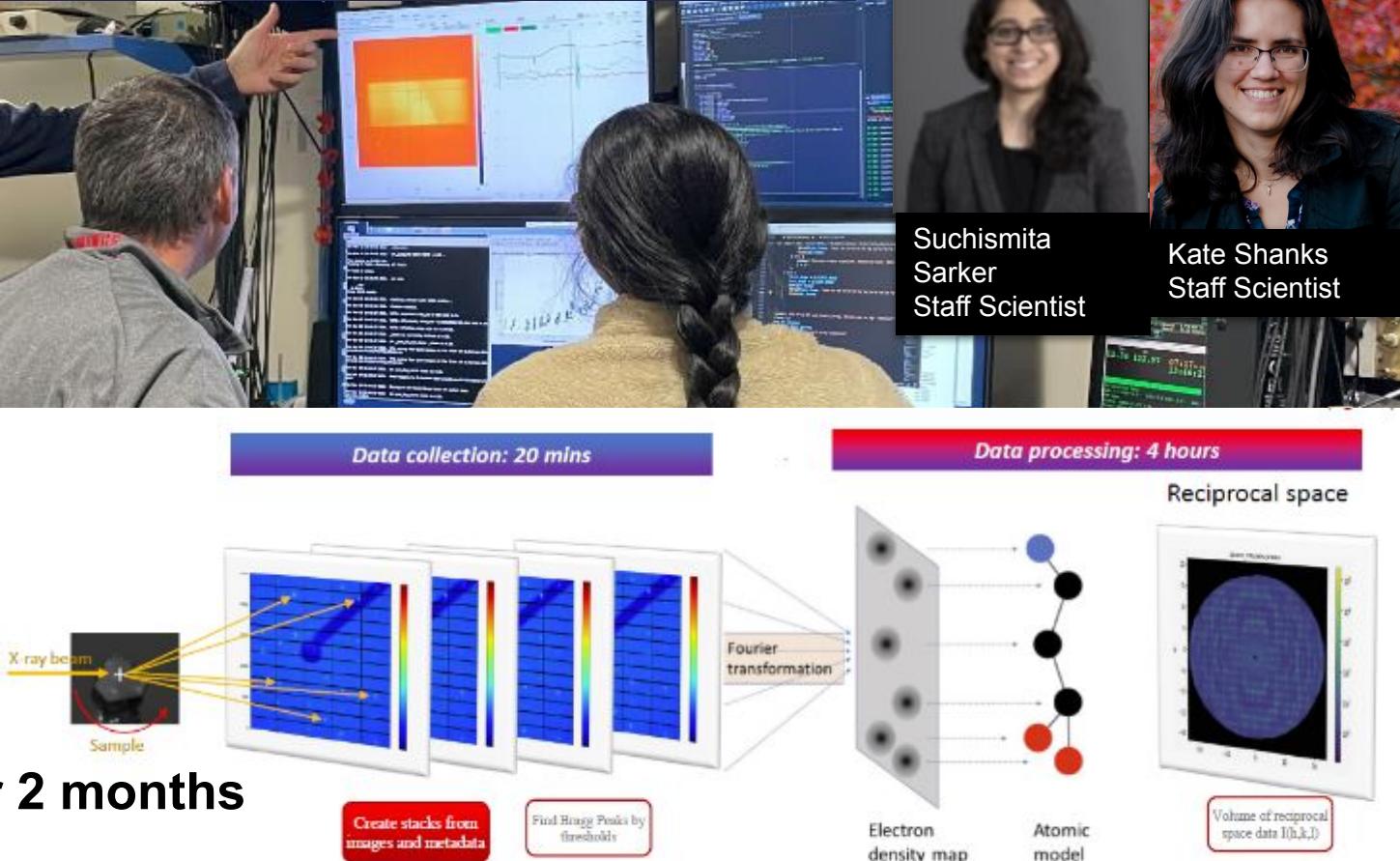


# Equity in the Use of Experimental Facilities

Establishing a comprehensive workflow from experimental facilities to the end-user data analysis.

- Cornell High Energy Synchrotron Source: Quantum Materials Beamline (6 of the 8 lines)
- Real-time data access and sharing
- Steering experiments
- AI-driven workflows
- Remote team collaboration
- Optimize effective use of scientists' time and national resource
- Publish data with no delay (e.g., Materials Commons) ☐ in real-time

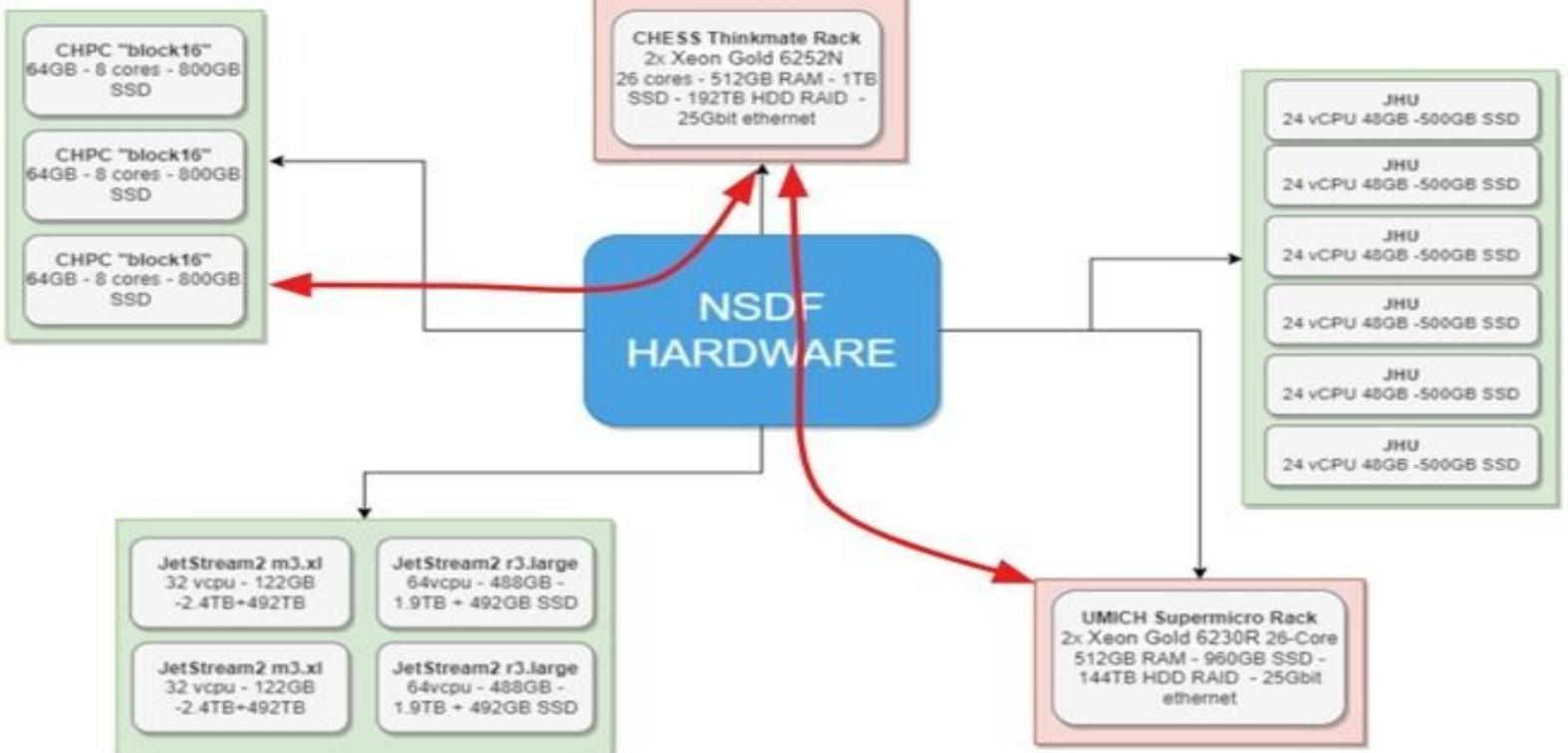
200TB data, 10M files in 2 beamlines over 2 months





Werner Sun · 2nd  
IT Director, Cornell Laboratory for Accelerator-Based Sciences and  
Education at Cornell University  
Ithaca, New York, United States · [Contact info](#)

# Direct Data Access & Steering from Cornell High Energy Synchrotron Source (CHESS)



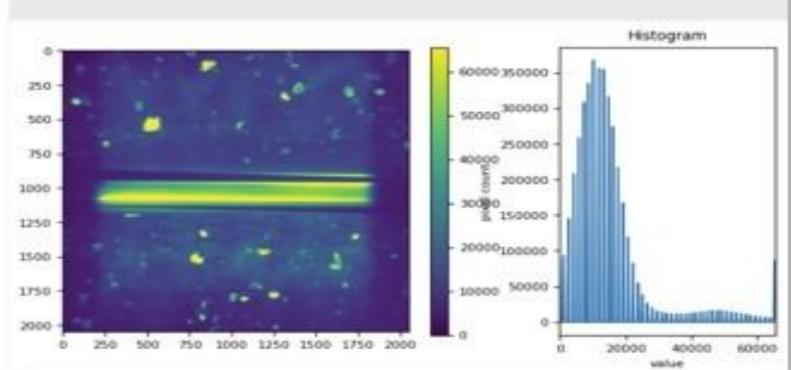
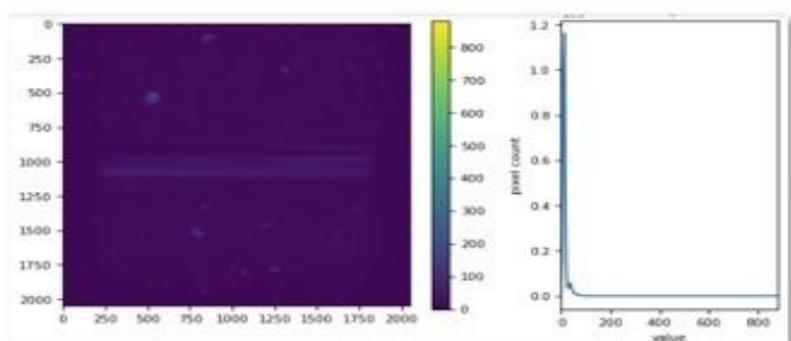
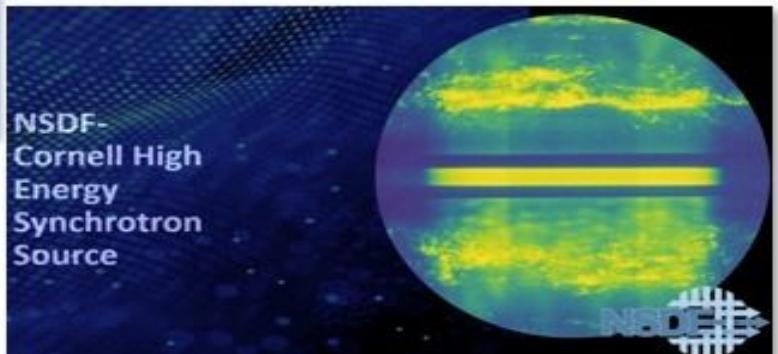
~1500 datasets

Uncompressed size 7.7TiB  
Compressed size 2.4TiB



CHESS

CORNELL HIGH ENERGY SYNCHROTRON SOURCE





The screenshot shows a mobile device displaying a web page from nsdf01.classe.cornell.edu. At the top, there's a header with the time (1:41), battery level (47%), and signal strength. Below the header, the URL is shown as nsdf01.classe.cornell... nsdf01.classe.cornell.edu. The main content area displays the "CHESS stats" dashboard, which includes a 3D wireframe visualization of a dataset and a table of data entries. The table has columns for id3a, id4b, and nsdf. The data entries include rows for various users like berman-3804-a, capolungo-3255-a, gopalan, greven-3798-a, lee-3365-c, pagari-3379-c, usn-3757-a, gregory-3844-a, gas-3858-a, ortiz-3822-a, and polihore-3828-s, each with links to a dashboard and a stat page.

A portrait of Suchismita (Suchi) Sarkar, a staff scientist at Cornell High Energy Synchrotron Source (CHESS). She is wearing glasses and a dark blazer over a black top. The background is a plain, light-colored wall. Overlaid on the image is a large blue rectangular box containing white text.

**“You can access data  
from anywhere on any  
devices.”**

— Suchismita (Suchi) Sarkar  
Staff Scientist, CHESS

**CHESS**  
Cornell High Energy Synchrotron Source

The screenshot shows a mobile device displaying a web page from nsdf01.classe.cornell.edu. The URL is ViSUS nsdf-group/far... nsdf01.classe.cornell.edu. The main content area shows a scientific visualization using the ViSUS software. It features a 3D plot of a complex, multi-layered structure, likely a crystal or a material sample, with a color scale ranging from purple to yellow. The plot is overlaid with numerous thin, vertical lines, possibly representing reflections or data points. The interface includes various controls and a color bar on the right side.



Kate Shanks, CHESS



STATUS SCIENCE

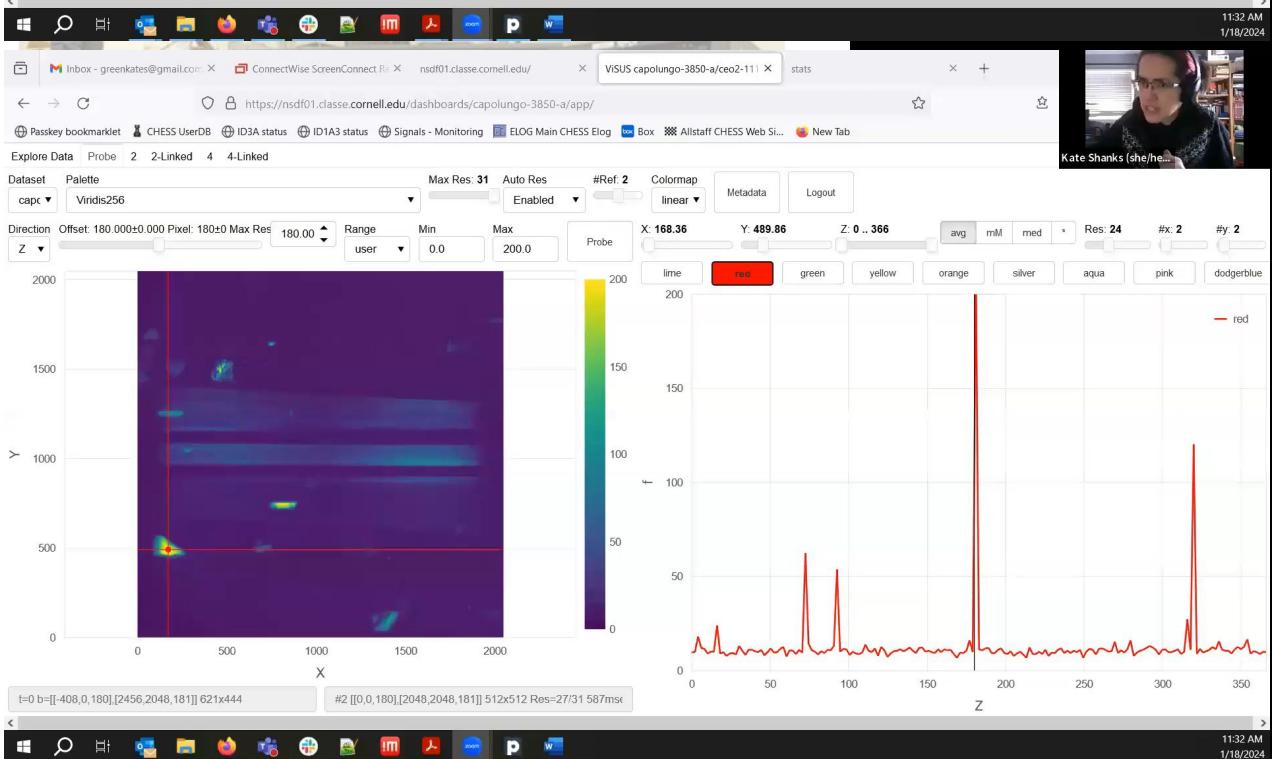
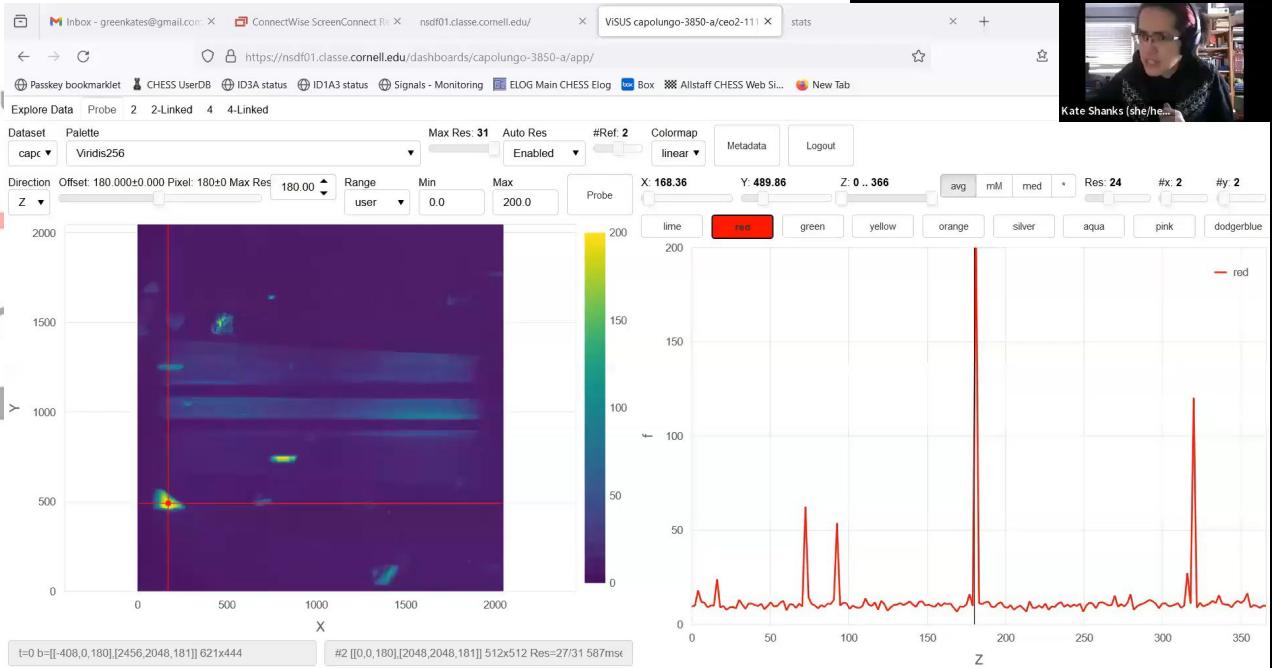
# National Science Data Fabric to Democratize Data-Driven

January 11, 2024 | Savan DeSouza

"It's much easier to pop open one of these image stacks on here, in your web browser, as opposed to trying to open it up on ImageJ. If you do that on the station computer, sometimes you crash the computer. So this decouples that and lowers the barrier to examining the datasets on the fly.."

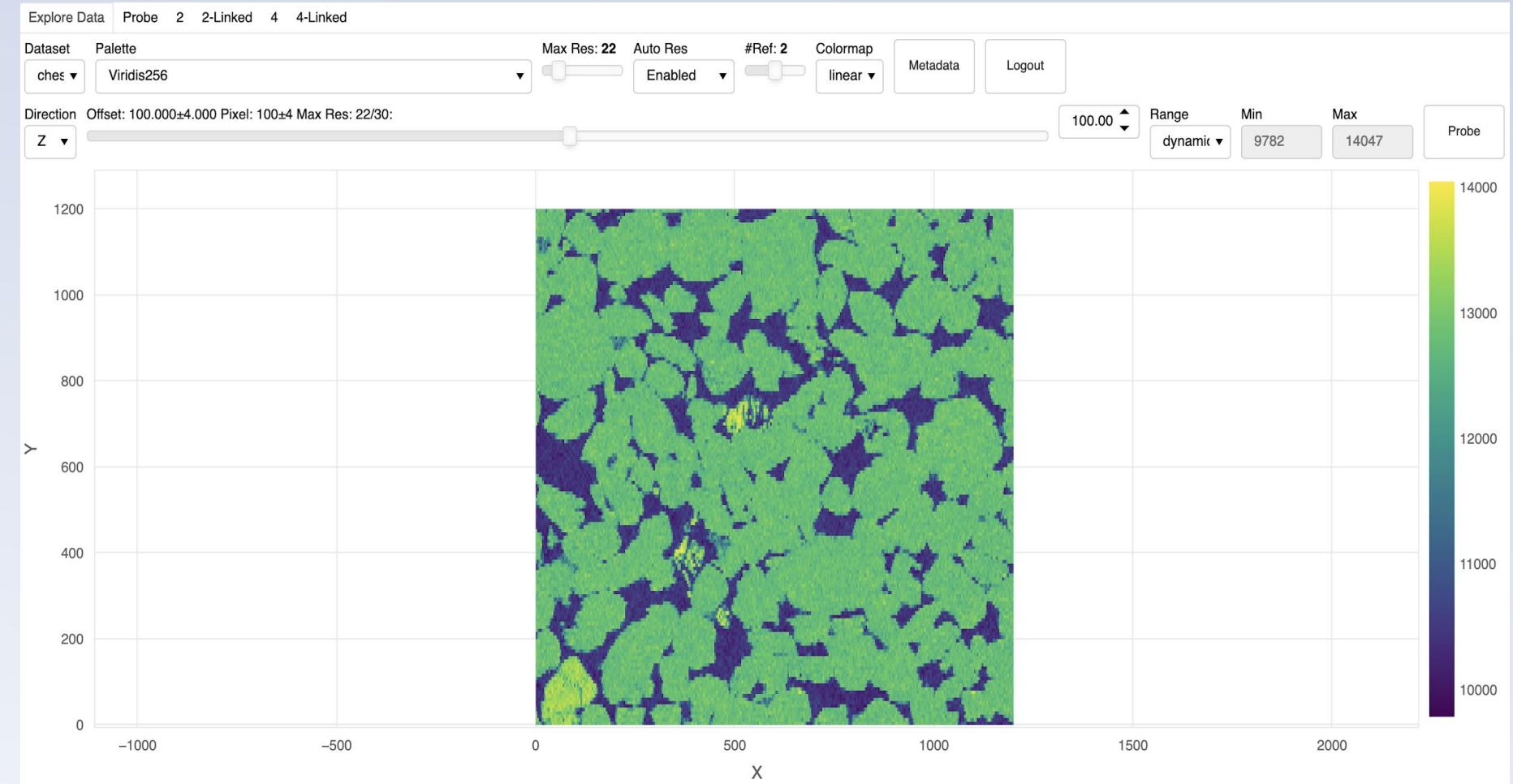
-Kate Shanks, FAST beamline scientist

<https://nation-d-cheze-data-repository>





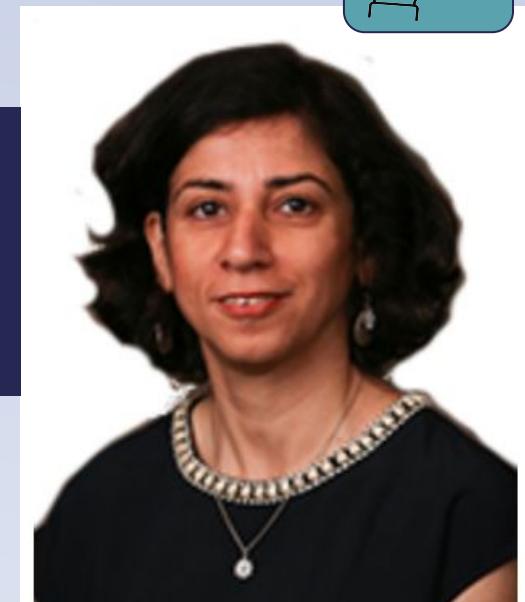
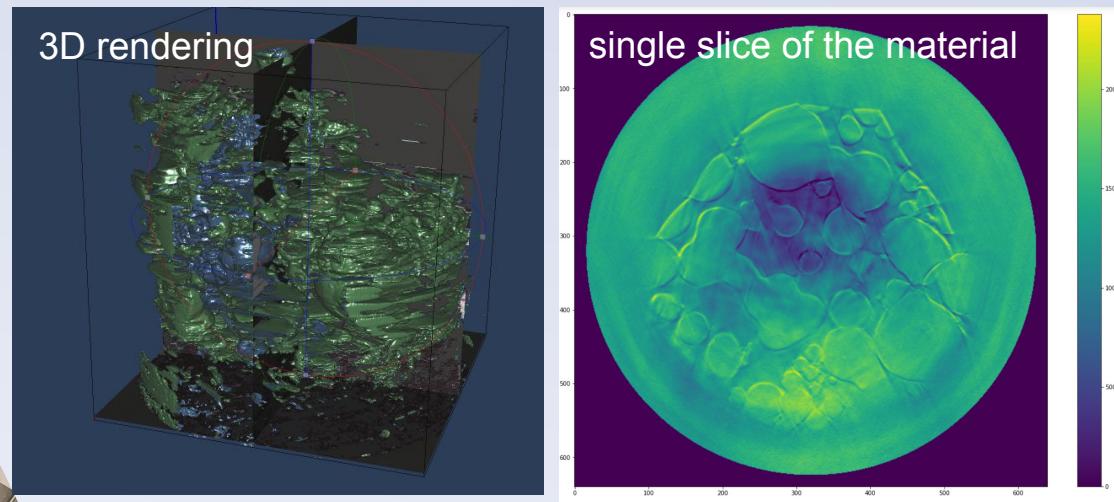
# Check Out the CHESS NSDF-Dashboard





# Equity in Research with National Labs

Facilitating rapid processing of 100+ terabytes of data by compact virtual laboratories, achieving in days what would take months for moving data between national labs and universities.

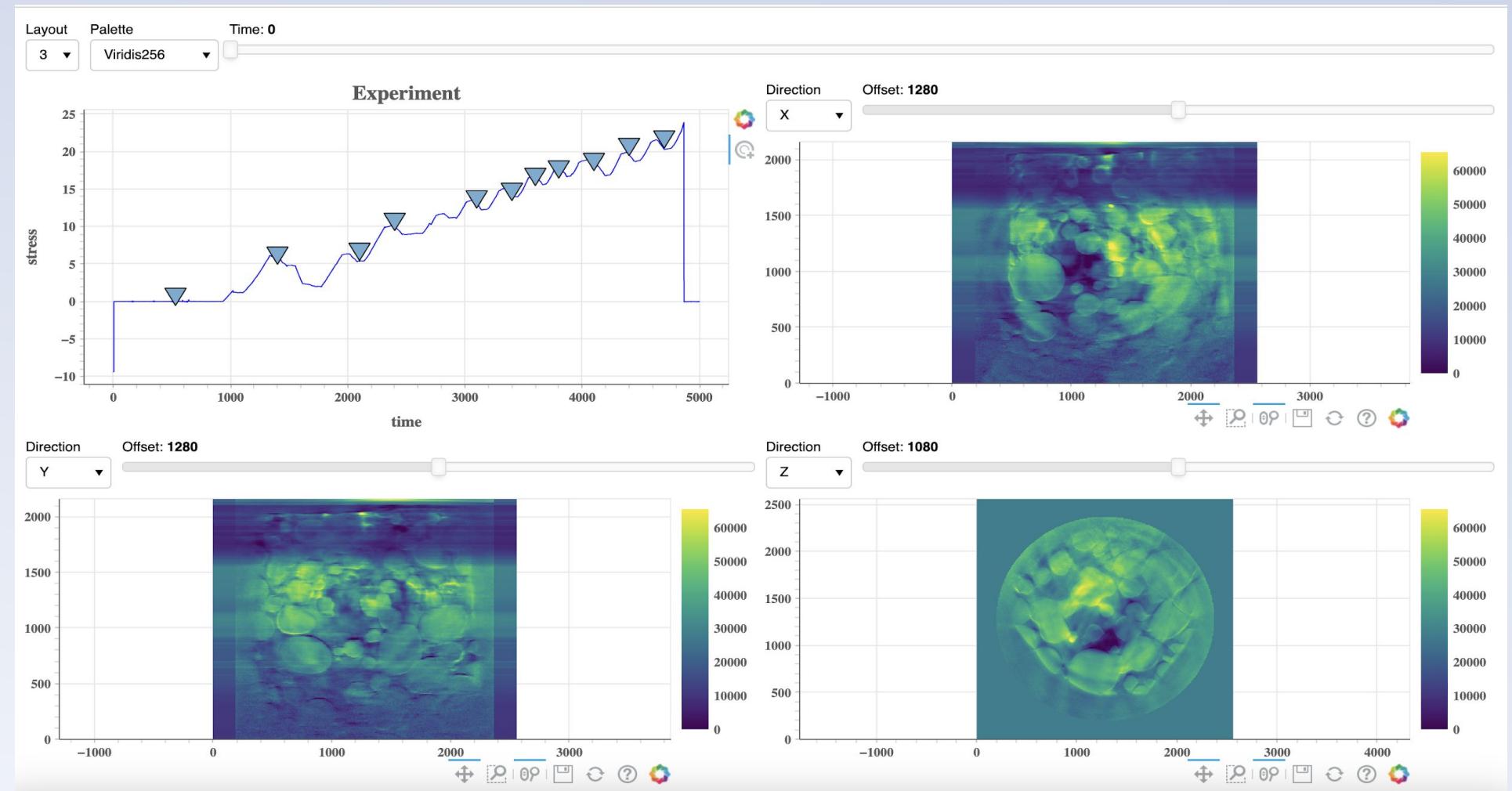


Pania Newell  
Professor in Mechanical  
Engineering at UoU

Analysis of **the internal structure of porous silica materials** at different stress levels  
Over **400TB of generate data** and more than 200 machines used on CloudLab, Chameleon, & AWS  
<http://services.nationalsciencedatafabric.org/materialscience>



# Check Out the Material Science NSDF-Dashboard





# Equity in Accessing Full Data in Real-Time from HPC and Cloud



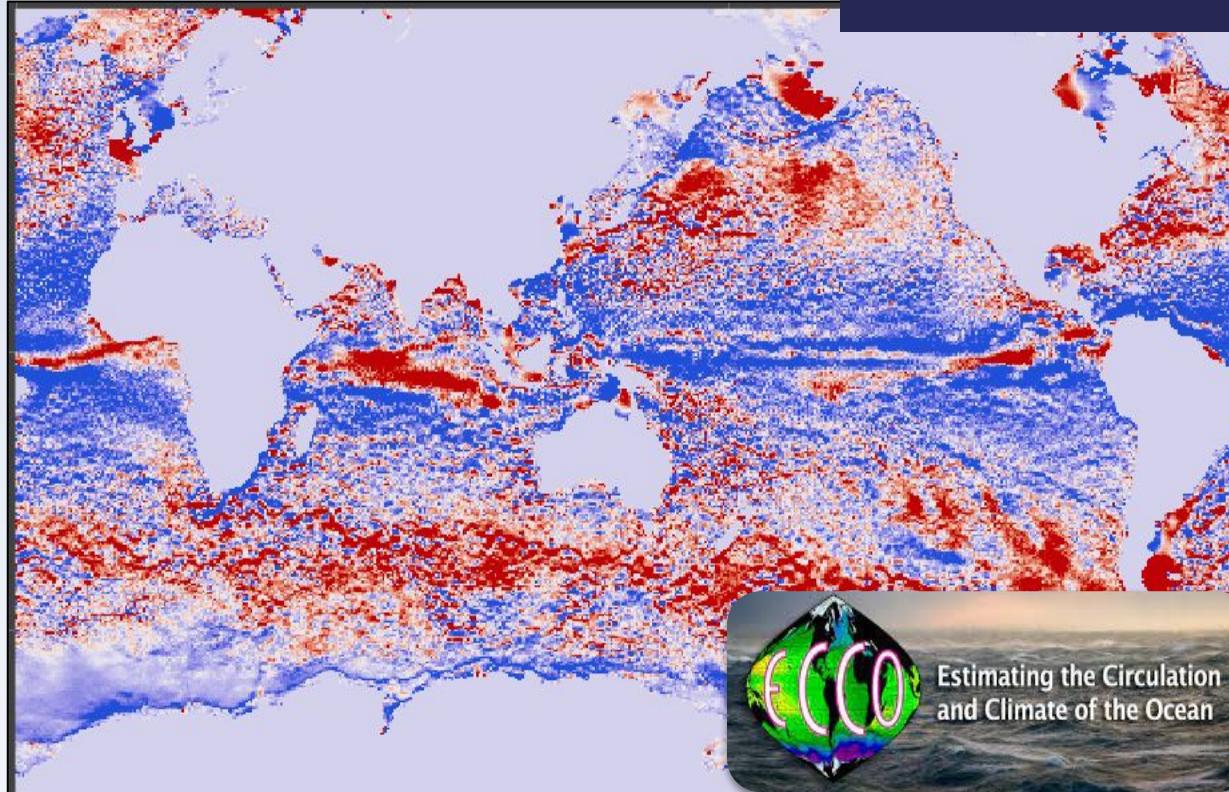
Nina McCurdy · 1st



NASA Ames Research Center

Data Analysis and Visualization Scientist at NASA Ames Research Center

Enabling streaming data access and processing in a distributed cyberinfrastructure for NASA climate modeling data (from Pleiades Supercomputer & Cloud)



LLC4320 2.8 PB Ocean dataset  
Model of ocean circulation and global ocean data to study the circulation role in climate changes

Requirements:

- Access the full data
- Overcome limitations in computational power
- Provide real-time processing capabilities.



National Science Data Fabric



www.sci.utah.edu



THE UNIVERSITY OF  
TENNESSEE  
KNOXVILLE



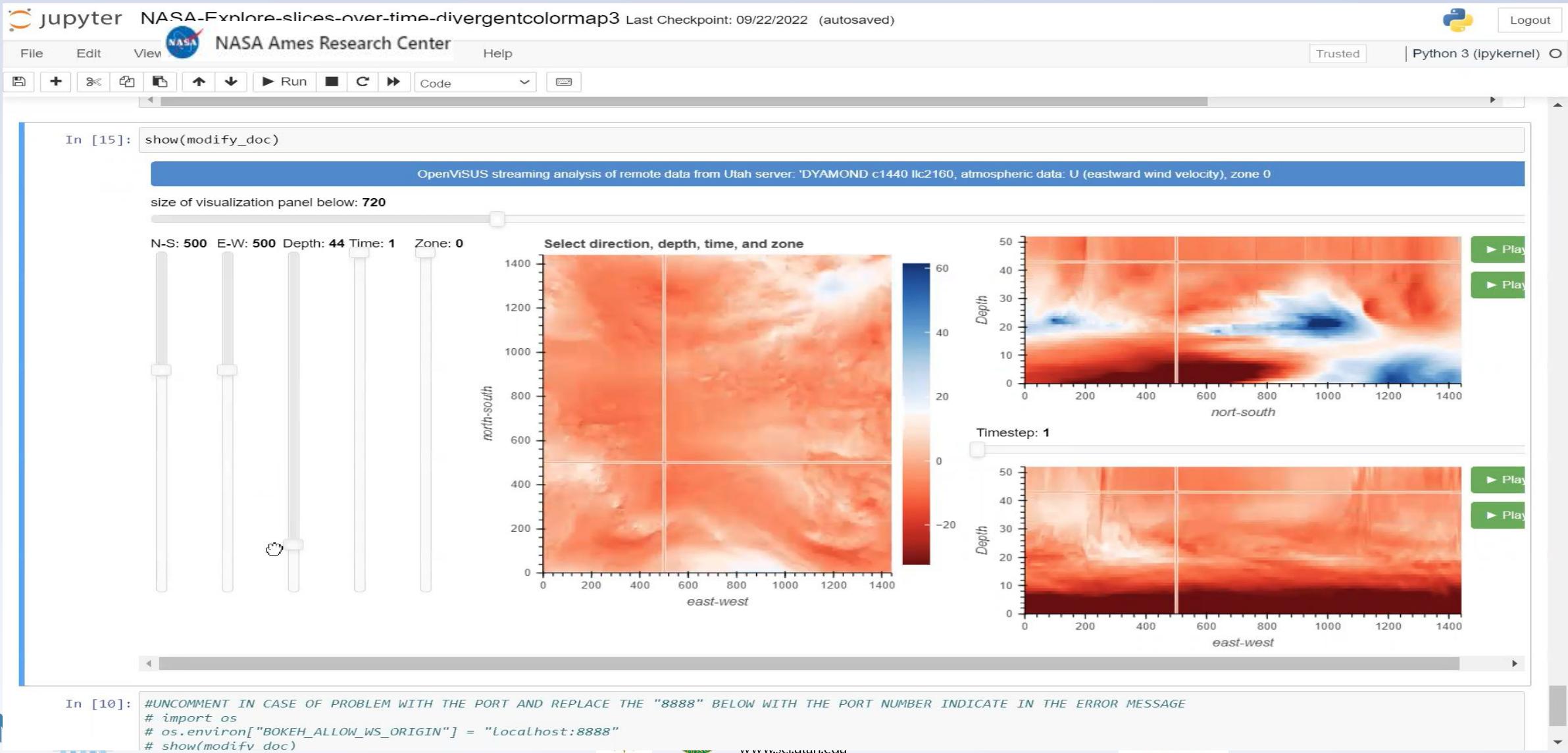
Powered by VISUS



JOHNS HOPKINS  
UNIVERSITY

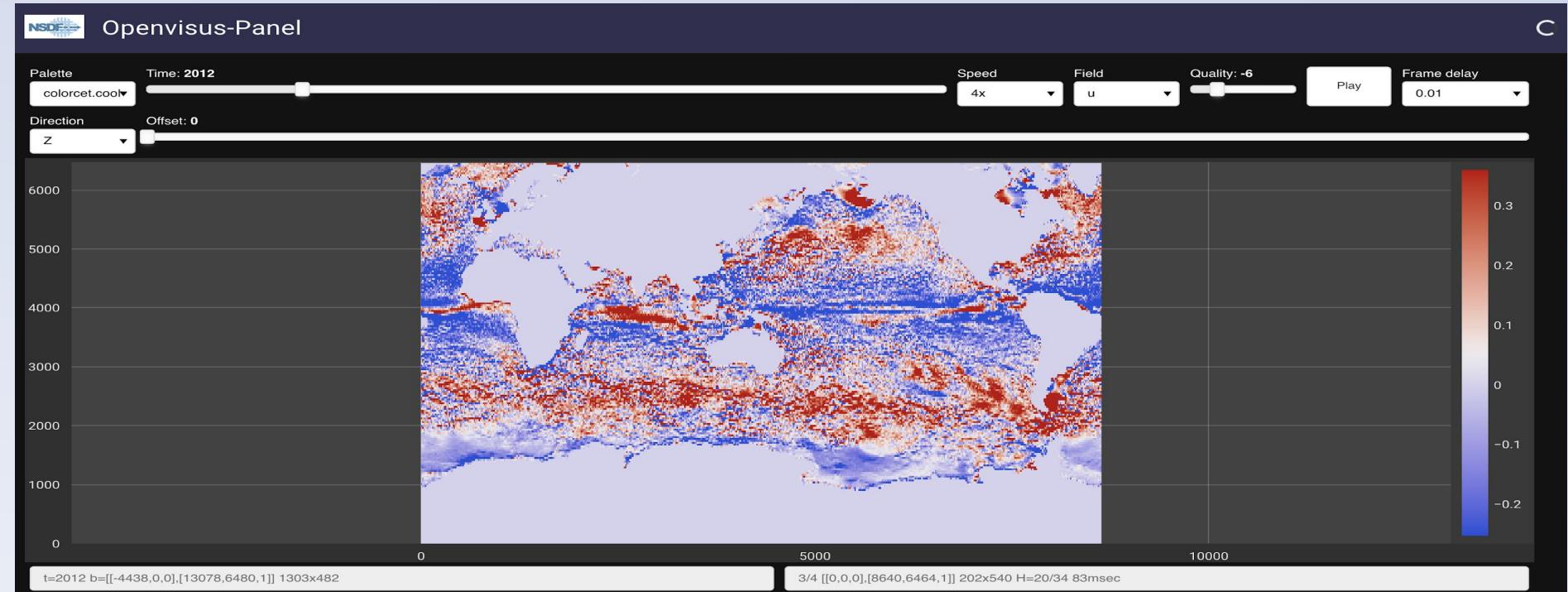


# Streaming Data Access for Run-Time Analysis





# Check Out NASA Ocean NSDF-Dashboard





# Session II: Training on using NSDF Services for End-to-End Analysis of Scientific Data

Democratizing Access and Use of Large-scale Data



National Science Data Fabric



[www.sci.utah.edu](http://www.sci.utah.edu)



THE UNIVERSITY OF  
TENNESSEE  
KNOXVILLE



Powered by ViSUS



SDSC

IBM





# Four-Step Workflow Tutorial

This tutorial showcases the capabilities of NSDF, guiding you through a **four-step modular workflow** that leverages OpenVisus services to analyze a geospatial dataset generated with GEOTiled.

## Step 1: Data Generation

Collect DEMs from the United States Geological Survey (USGS). Process them with GEOTiled or upload the data from public or private storage.

## Step 2: Conversion to IDX Format

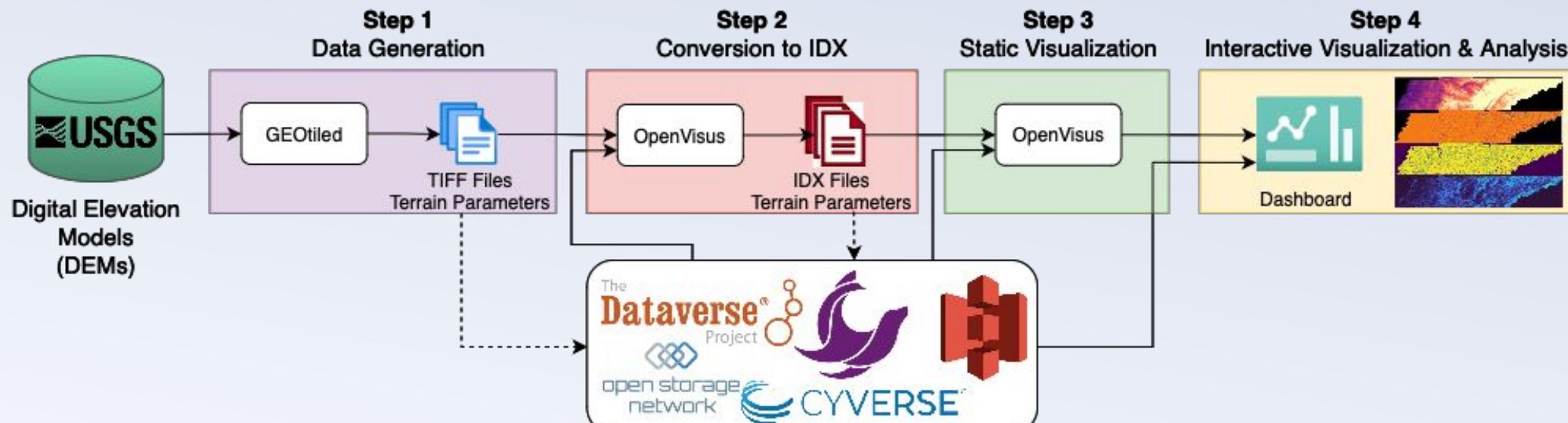
Convert files from TIFF to IDX (the format used by OpenVisus), preserving accuracy but reducing size. Store IDX files in public or private storage.

## Step 3: Static Visualization

Statically visualize the terrain parameters in OpenVisus. Validate accuracy of IDX-based images with the TIFF-based images.

## Step 4: Interactive Visualization & Analysis

Launch dashboard for interacting with large-scale data to access subregions of the original dataset for ad hoc analysis



(1) Your codespace is fully loaded. Now you should see this screen



Enabling Scientific Discovery: Harnessing the Power of the National Science Data Fabric for Large-Scale Data Analysis

**Abstract**

In this interactive half-day tutorial, participants explore the advanced applications of the National Science Data Fabric (NSDF) services and comprehensive strategies for end-to-end scientific data analysis. The tutorial targets a broad audience, from researchers and students to developers and scientists, each finding valuable insights into managing and analyzing large datasets, with a particular focus on datasets exceeding 100TB. Attendees gain hands-on experience constructing modular workflows, leveraging public and private data storage and streaming solutions, and deploying sophisticated visualization and analysis dashboards for scientific discovery. The tutorial highlights NSDF's role in supporting the VIS conference's themes by providing scalable solutions for advances in visualization and visual analytics. It covers various topics, from an overview of NSDF's capabilities to addressing common pain points in data analysis to intermediate hands-on exercises using NSDF services for Earth science data and advanced applications, including handling and visualizing massive datasets in domains requiring high-resolution data management. Participants leave a deeper understanding of how NSDF services integrate into their research workflows to enhance data accessibility, sharing, and collaborative scientific discovery. This tutorial advances the knowledge of data-intensive computing and empowers attendees to harness the full potential of NSDF in their fields.

**Agenda**

| Session  | Duration | Objective     |
|----------|----------|---------------|
| PROBLEMS | OUTPUT   | DEBUG CONSOLE |
| TERMINAL | PORTS    | COMMENTS      |

(base) root@codespaces-3c1c53:/workspaces/Tutorial\_2024\_IEEE\_VIS#

(3) Look for the *Preparing your Environment* section and press the play button



EXPLORER

- ACM\_SUMMER SCHOOL\_20...
- .devcontainer
- .github/workflows
- ! main.yml
- hands-on
  - session II
    - Materials
    - 1.Tutorial.ipynb
    - 2.Explore\_Data.ipynb
  - session III
  - slides
  - videos
- .gitignore
- README.md

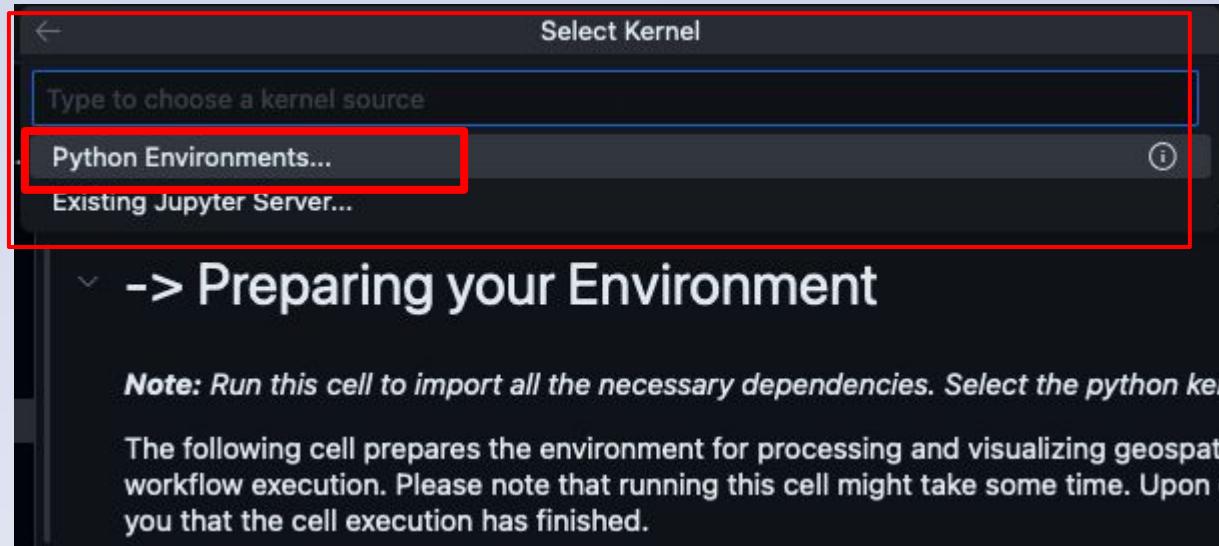
```
import geotiled as gt
from pathlib import Path
import glob
import os
import shutil
import multiprocessing
import OpenVisus as ov
import numpy as np
import requests
import json
from matplotlib import pyplot as plt
from tqdm import tqdm

# To silence a deprecation warning.
gt.gdal.UseExceptions()
```

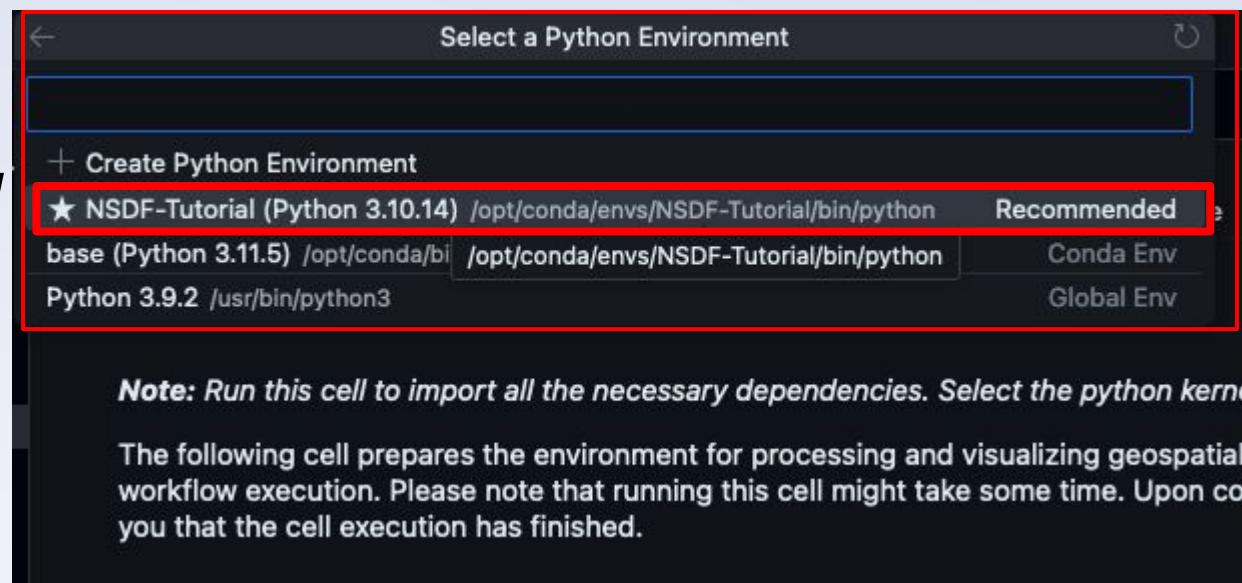
Instructions continue



(4) A list to →  
Select Kernel  
should pop up.  
Select *Python Environments...*



(5) Select the  
★NSDF-Tutorial  
option →



(6) After running the  
*Preparing your  
Environment* cell,  
you should see a  
message saying it  
was successfully  
prepared

A screenshot of a Jupyter Notebook cell showing the output of the code: '# You have successfully prepared your environment.' followed by 'print("You have successfully prepared your environment.")'. The output shows the message 'You have successfully prepared your environment.' and a duration of '1.2s'.



# Step 1: Data Generation with GEOTiled

## Step 1: Data Generation

Collect DEMs from the United States Geological Survey (USGS) and process them with GEOTiled or upload the data from public or private storage.

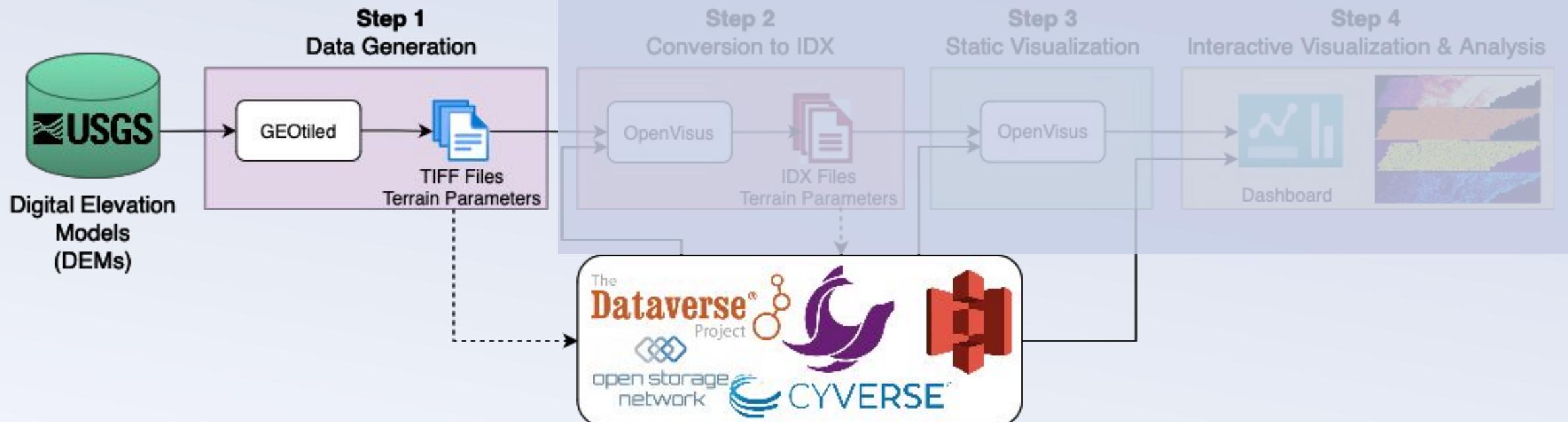
Step 1 provides **two options** to obtain data and generate the TIFF files before proceeding with Step 2

### Option A

Generating Data Using the SOMOSPIE Application Module

### Option B

Accessing data from Dataverse public commons



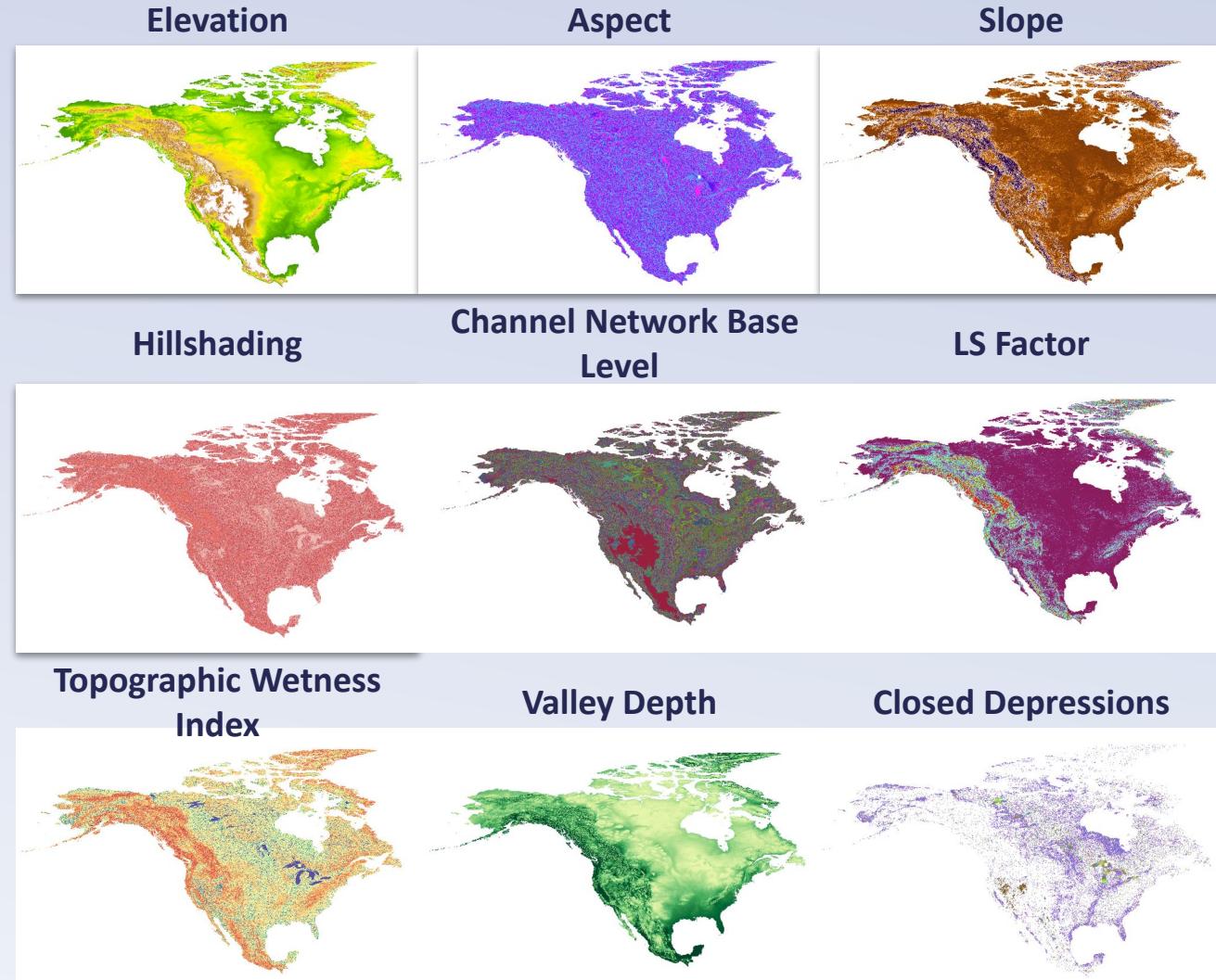


# Step 1: What are Terrain Parameters?

Terrain parameters (e.g., slope, aspect, hillshading, etc.) are **descriptions of surface form derived from Digital Elevation Models (DEM)**.

They play a **fundamental role** in applications such as **precision forestry and agriculture, and hydrology for landscape ecology**.

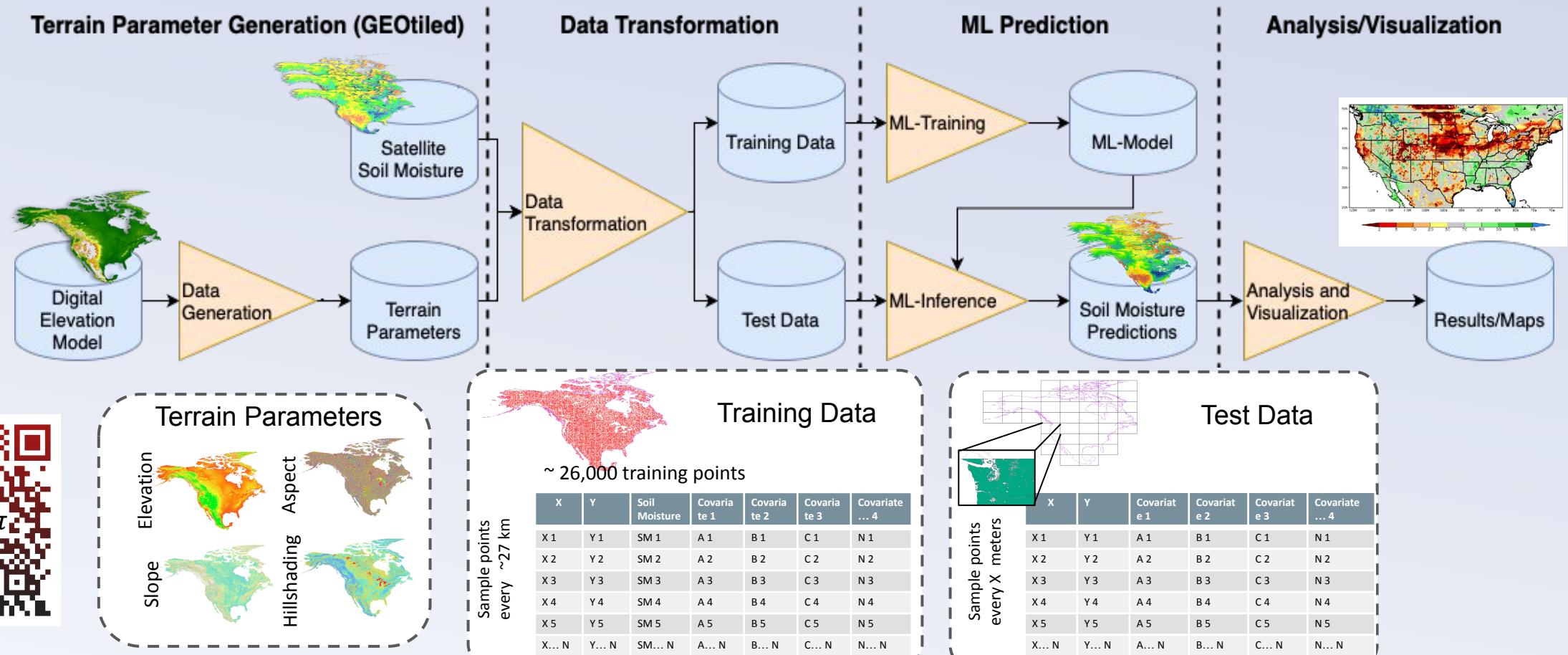
Generating **terrain parameters at high-resolution** is computationally **expensive**, hindering their accessibility by the scientific community





# Step 1: *SOMOS* $\pi$ Components

SOMOSPIE (SOil MOisture SPatial Inference Engine) has **four components** that empower scientists to generate, predict, and analyze **high-resolution topographic data**



National Science Data Fabric



www.sci.utah.edu



TENNESSEE

KNOXVILLE

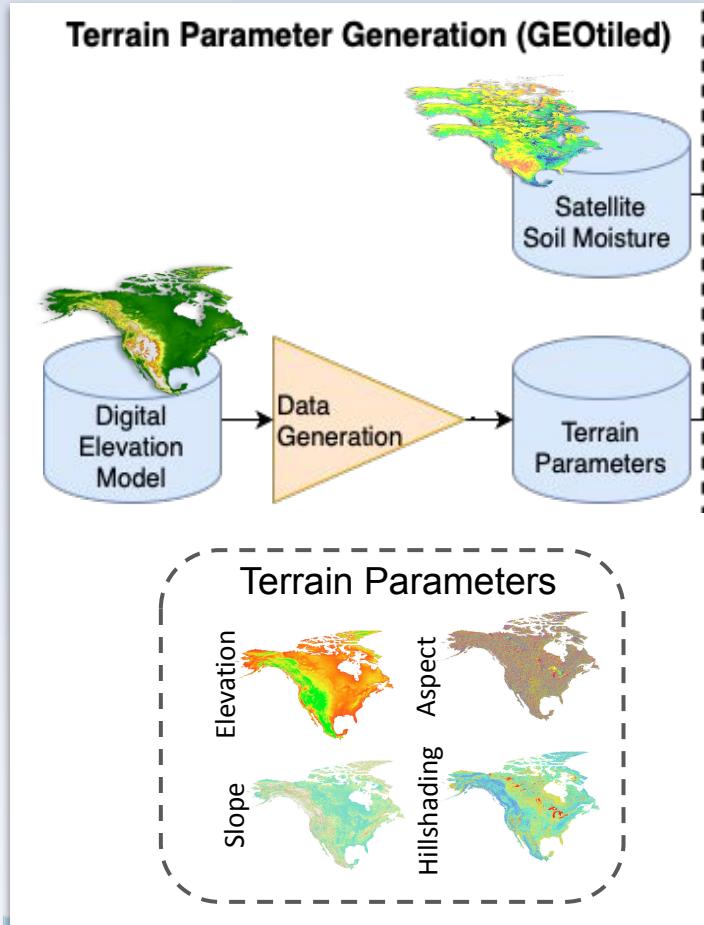


46

# Step 1: GE<sup>O</sup>tiled Terrain Generation

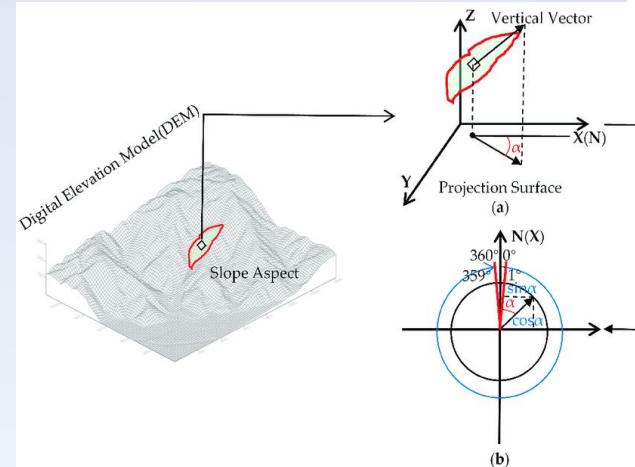


We expand on the first component, **GEOtiled**, that **computes high-resolution terrain parameters** using Digital Elevation Models (DEMs)

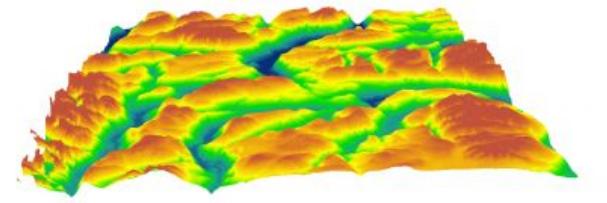


**GEOtiled** leverages data partitioning to accelerate the computation of terrain parameters from DEMs while preserving accuracy

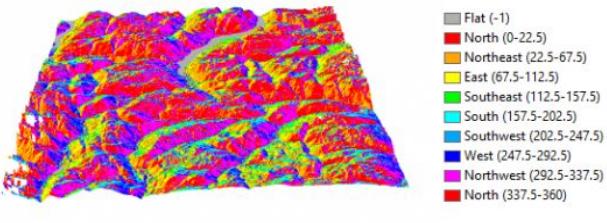
Computing High-resolution  
Terrain Parameters



Digital Elevation Model (DEM)



Terrain Parameter 1: Aspect



National Science Data Fabric



SDSC

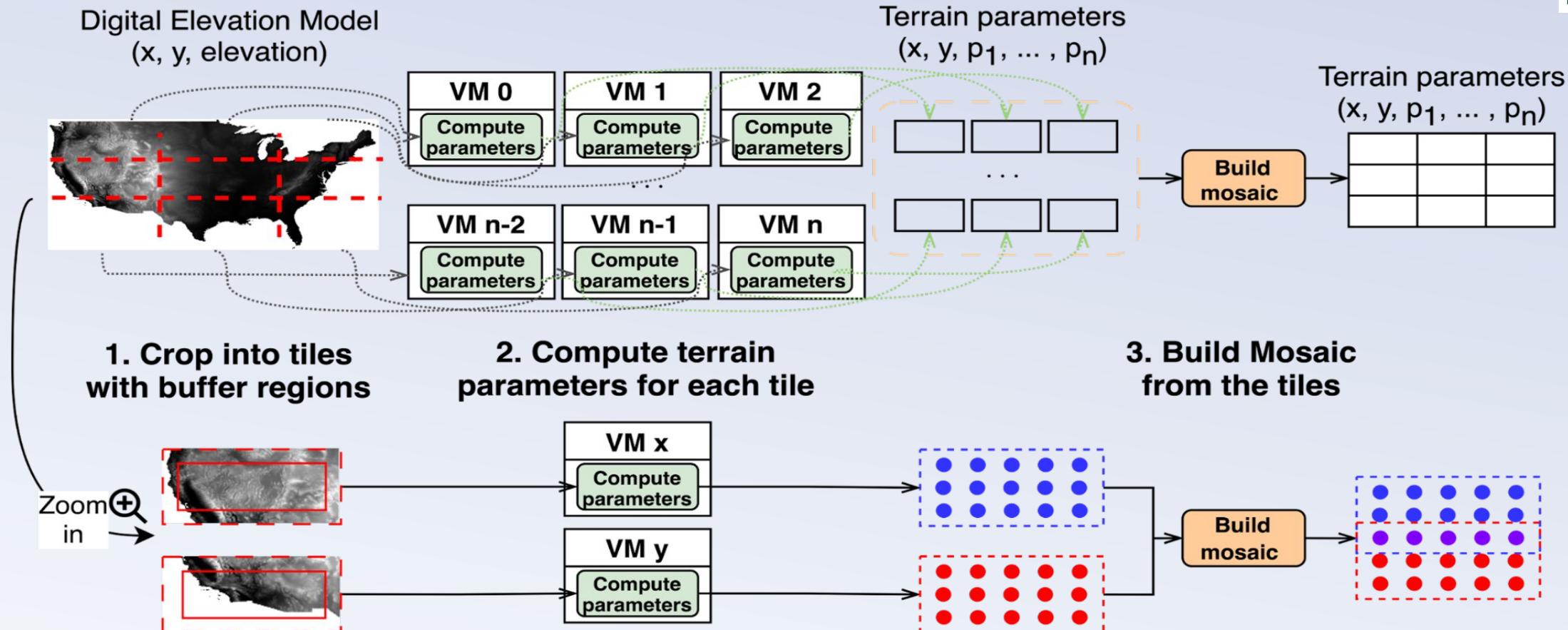
IBM





# Step 1: GEotiled Terrain Generation

GEotiled generates high-resolution terrain parameters at a large scale from a DEM. It has three stages: **(1)** we **crop DEM into tiles**, with buffer regions; **(2)** we **compute the terrain parameters** for each tile; **(3)** we **build a mosaic** from the tiles.



National Science Data Fabric



THE UNIVERSITY OF  
TENNESSEE  
KNOXVILLE



SDSC

IBM



# Step 1: Data Generation with GEOTiled



## Step 1: Data Generation

Collect DEMs from the United States Geological Survey (USGS) and process them with GEOTiled or upload the data from public or private storage.

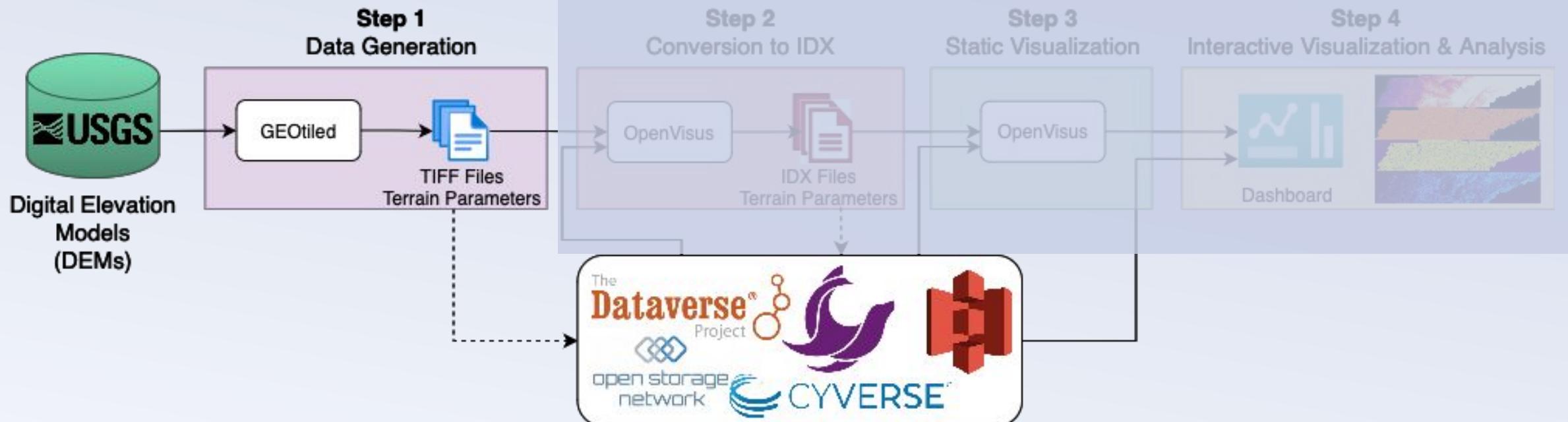
Step 1 provides **two options** to obtain data and generate the TIFF files before proceeding with Step 2

### Option A (~10 mins)

Generating Data Using the SOMOSPIE Application Module

### Option B (~3 mins)

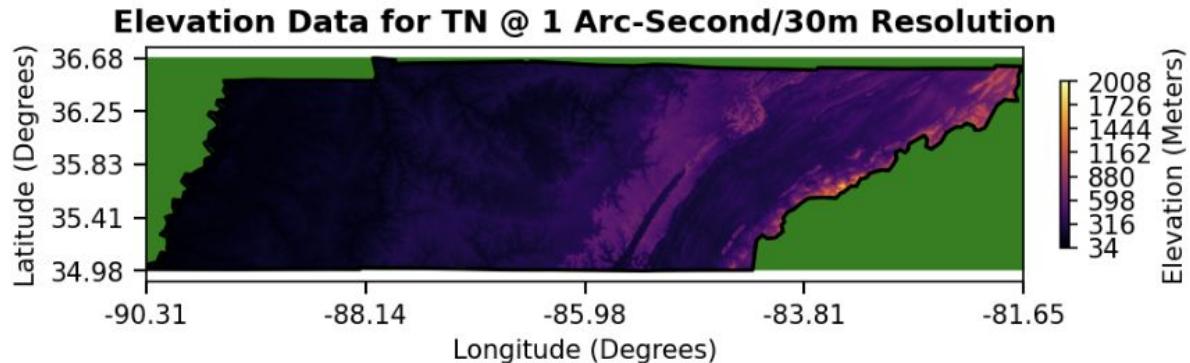
Accessing data from Dataverse public commons



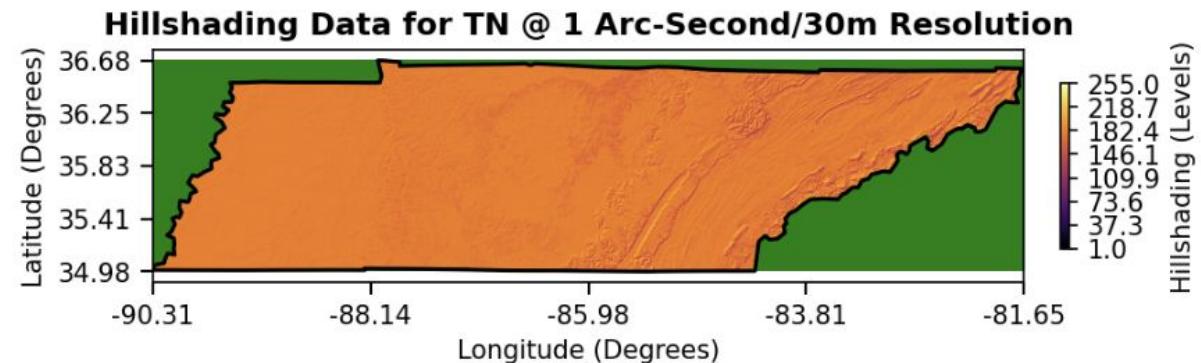
# Step 1: Data Generation with GEOTiled



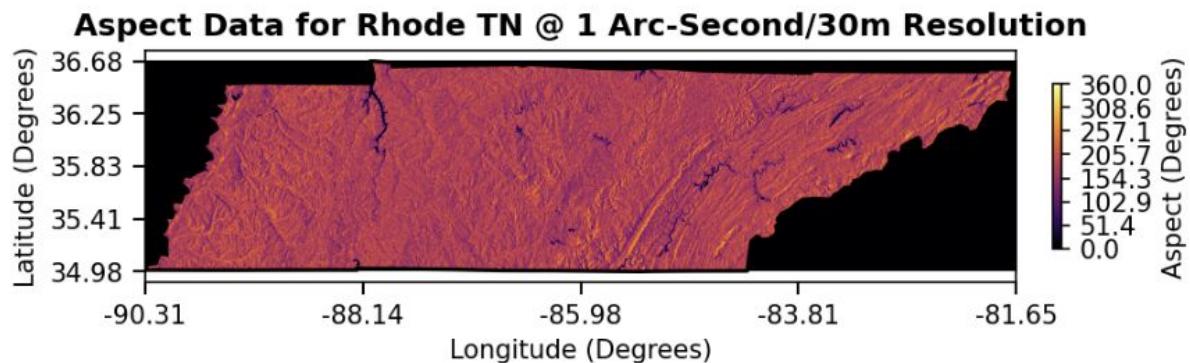
[Tutorial](#)



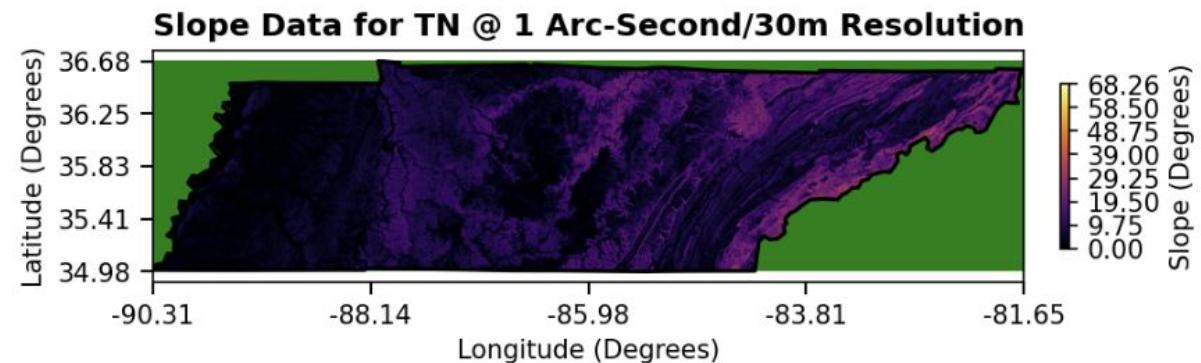
(a) Elevation - Terrain Parameter



(b) Hillshading - Terrain Parameter



(c) Aspect - Terrain Parameter



(d) Slope - Terrain Parameter



National Science Data Fabric



SDSC

IBM



# Step 2: Conversion to IDX



[Tutorial](#)

**OpenVisus** is a progressive cache-oblivious framework for large-scale data visualization

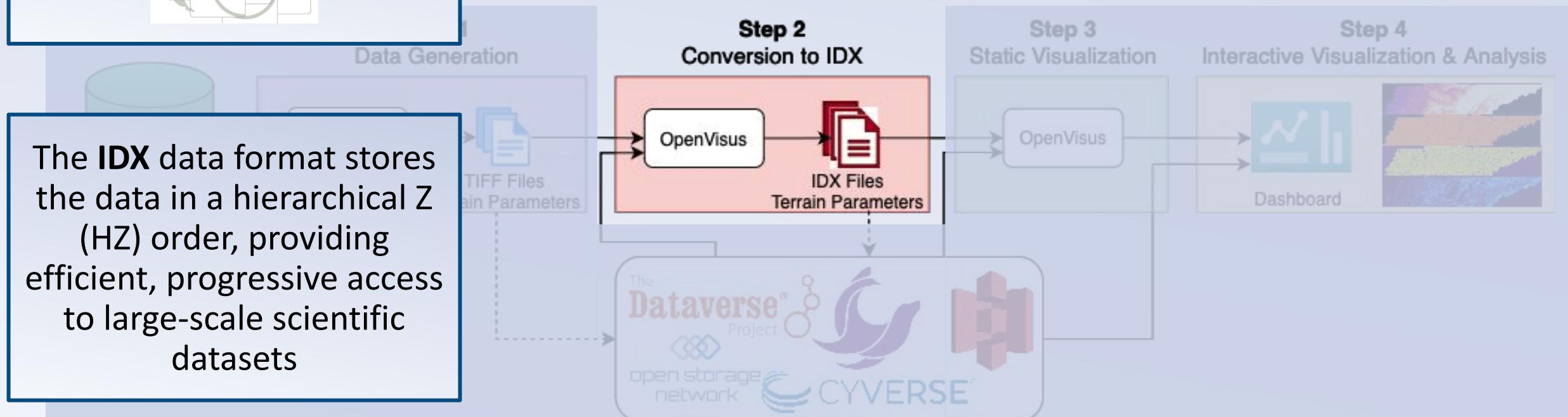


The **IDX** data format stores the data in a hierarchical Z (HZ) order, providing efficient, progressive access to large-scale scientific datasets

## Step 2: Conversion to IDX Format

Convert files from TIFF to IDX (the format used by **OpenVisus**), preserving accuracy but reducing size. Store **IDX** files in public or private storage.

Converting to **IDX** from TIFF format **reduces file size by 20%** while preserving accuracy



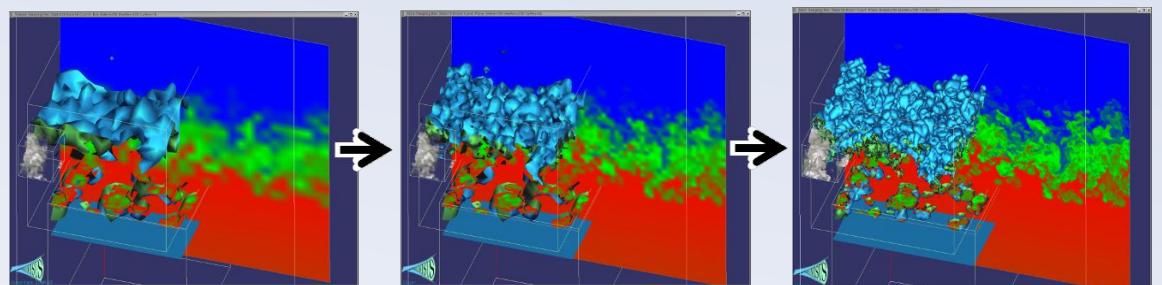


# Step 2: IDX Data Format

## Why IDX?

- The IDX data format provides **efficient, cache-oblivious, and progressive access** to large-scale scientific datasets.
- Data stored in IDX format can be visualized in an **interactive environment** allowing for meaningful explorations with **minimal resources**.
- IDX provides **scalability** across a wide range of running conditions like personal computers to distributed systems.

- Conversion to IDX is **not limited** to TIFF; it will work on other data formats like **NetCDF, HDF5, RGB, raw/binary**, and so on.
- IDX supports industry-standard lossless and lossy compression algorithms such as `zlib`, `zfp`, `lz4`.





# Step 3: Static Visualization

Step 3 provides **two options** to obtain data and collect the IDX files

## Option A

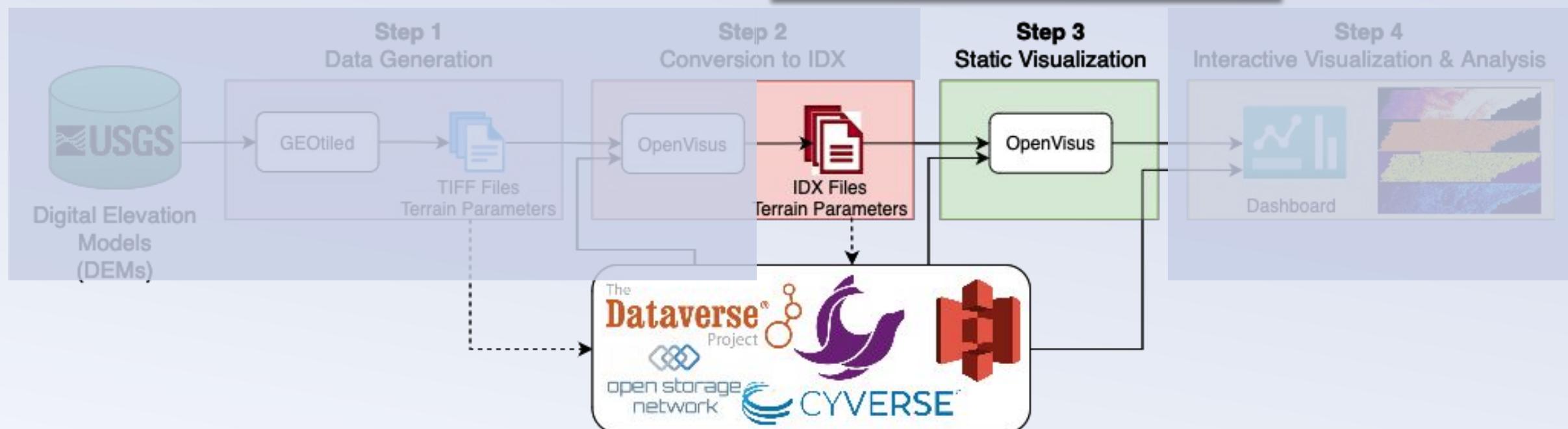
From local storage

## Option B

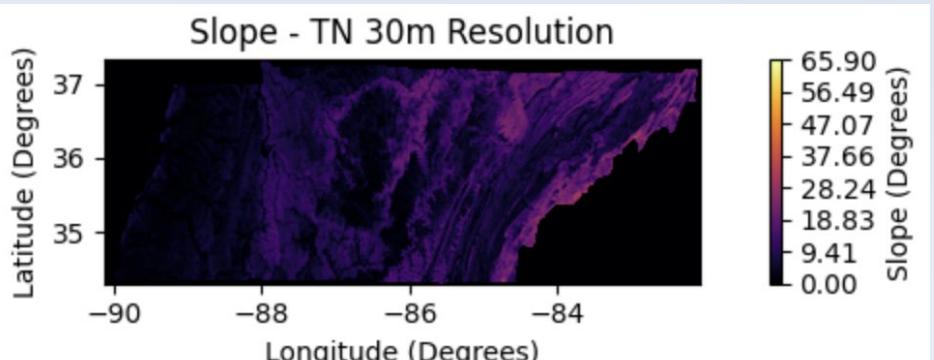
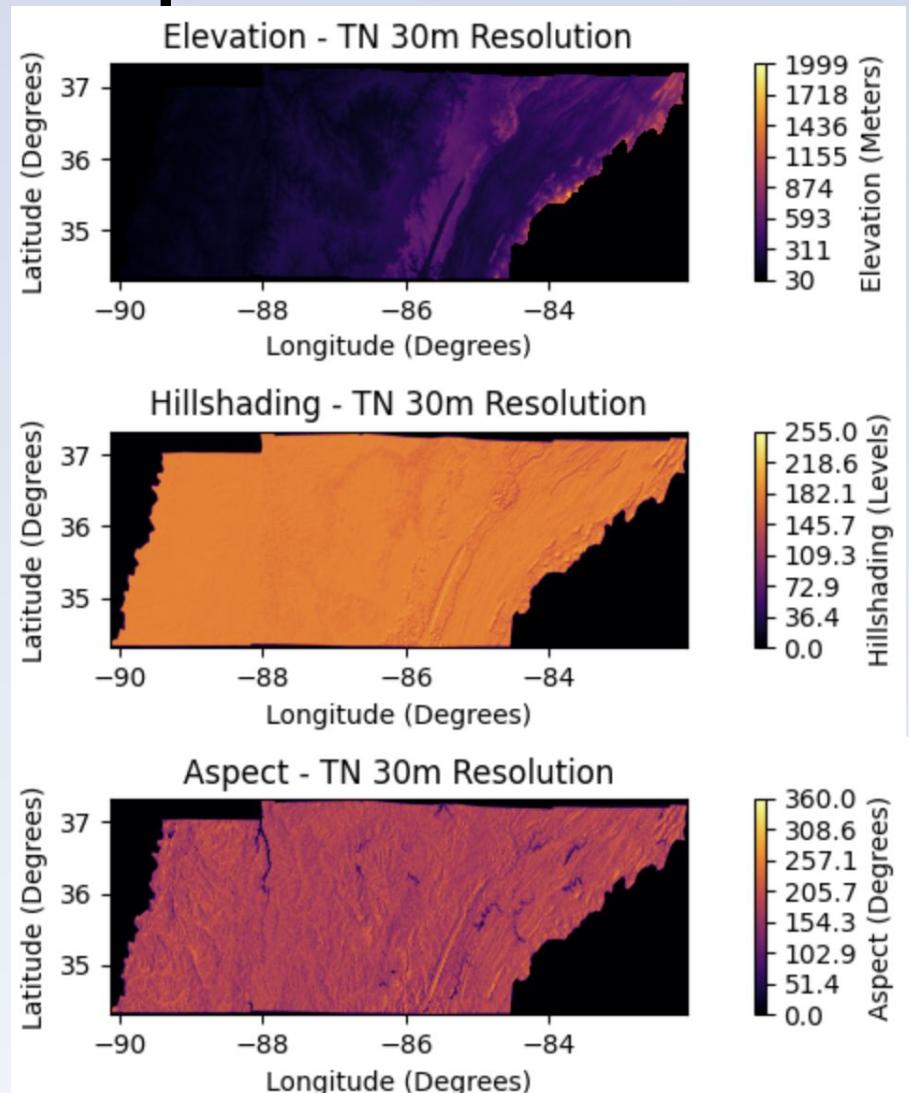
From Seal Storage

## Step 3: Static Visualization

Statically visualize the terrain parameters in OpenVisus. Validate the accuracy of IDX-based images with TIFF-based images.

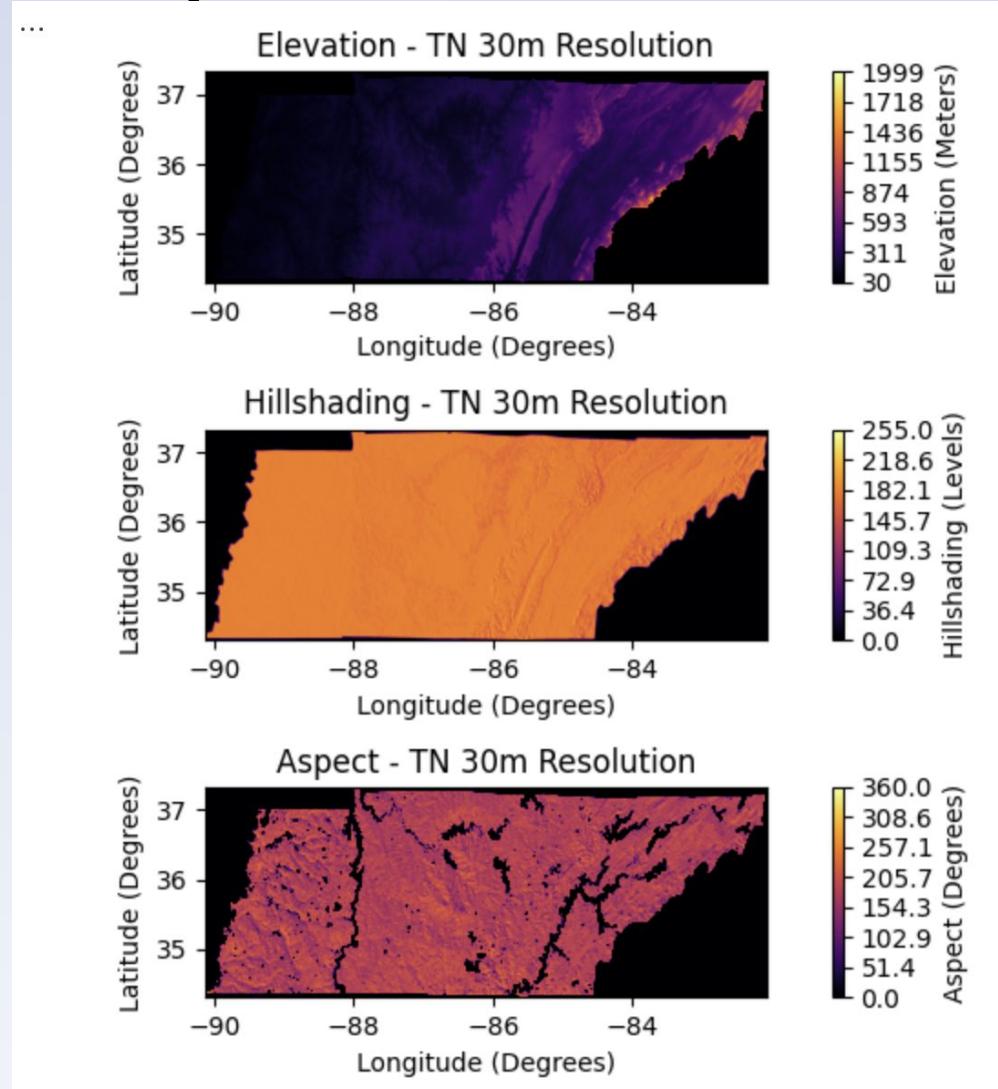


# Step 3: Static Visualization

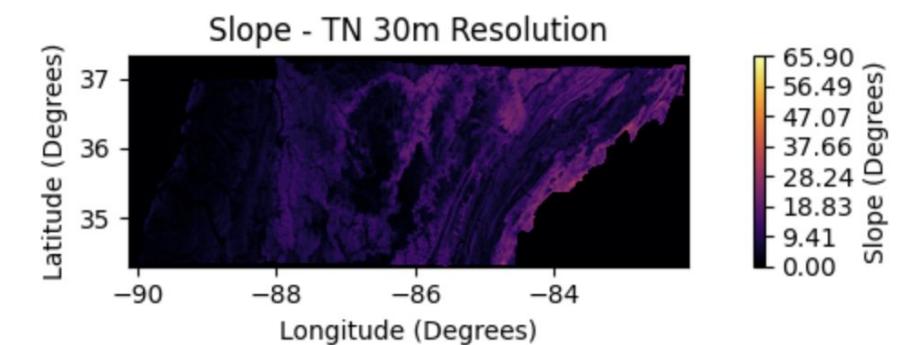


**Option A**  
From local storage

# Step 3: Static Visualization



**Option B**  
 From Seal Storage



... You have successfully visualized the IDX files.

# Step 4: Interactive Visualization & Analysis

Remotely **access** large datasets, **zoom** into specific areas, **select** and **crop** subregions of interest, **save** data locally in a Python-compatible format, and **analyze** the data for scientific discovery.

Step 4 provides **two options** to obtain data and collect the IDX files

## Option A

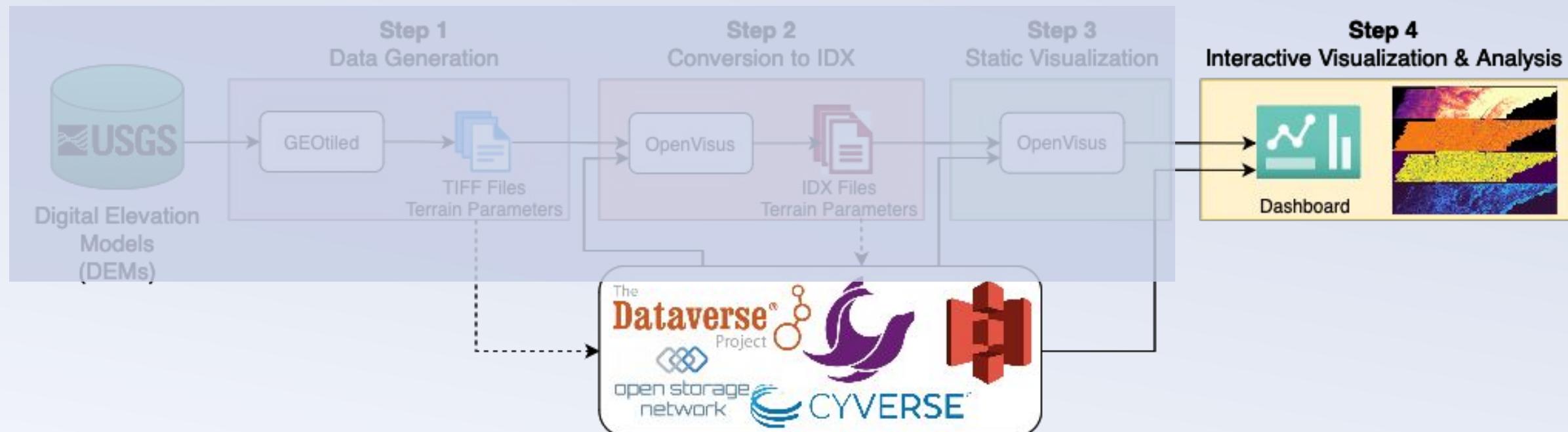
From local storage

## Option B

From Seal Storage

## Step 4: Interactive Visualization & Analysis

Launch dashboard for interacting with large-scale data to access subregions of the original dataset for ad hoc analysis.

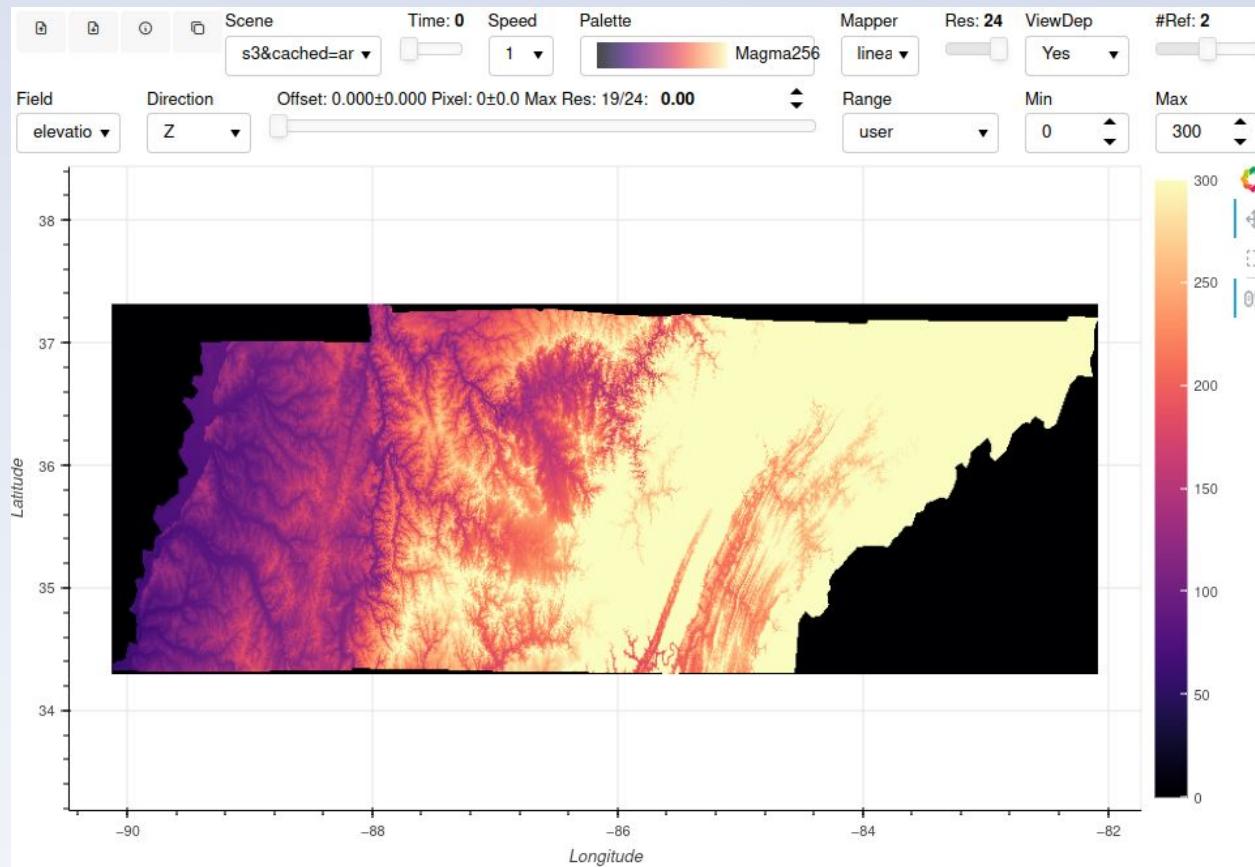




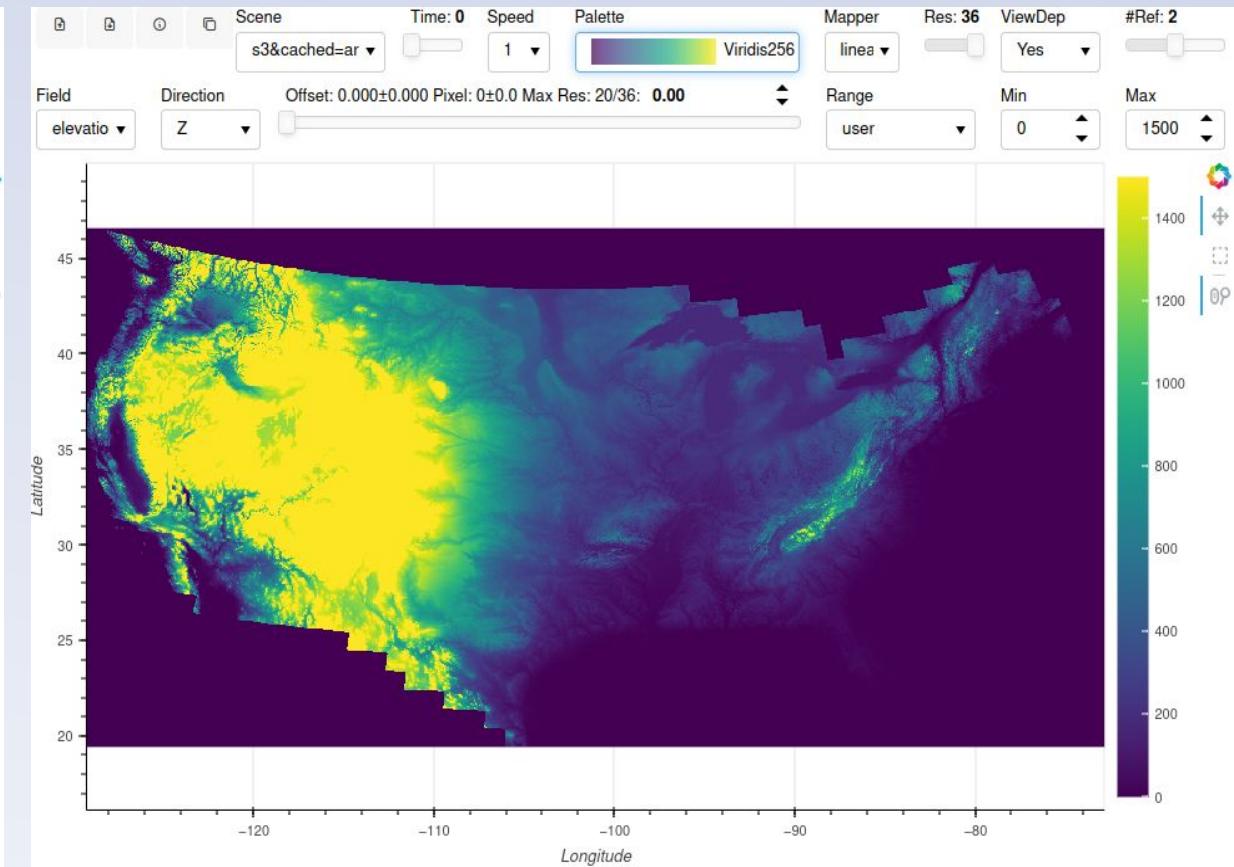
# Step 4: Geographical Regions

Visualize and analyze two geographical regions at 30 m resolution

State of Tennessee - 200 MB



Contiguous United States (CONUS) - 200 GB



# Step 4: Interactive Visualization & Analysis

Remotely **access** large datasets, **zoom** into specific areas, **select** and **crop** subregions of interest, **save** data locally in a Python-compatible format, and **analyze** the data for scientific discovery.

Step 4 provides **two options** to obtain data and collect the IDX files

## Option A

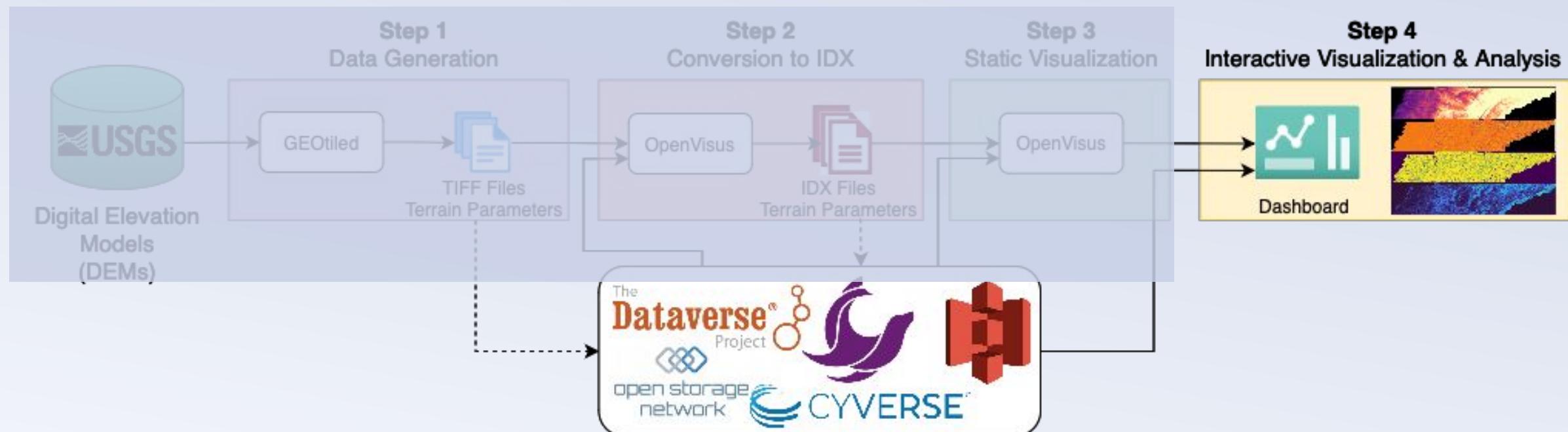
From local storage

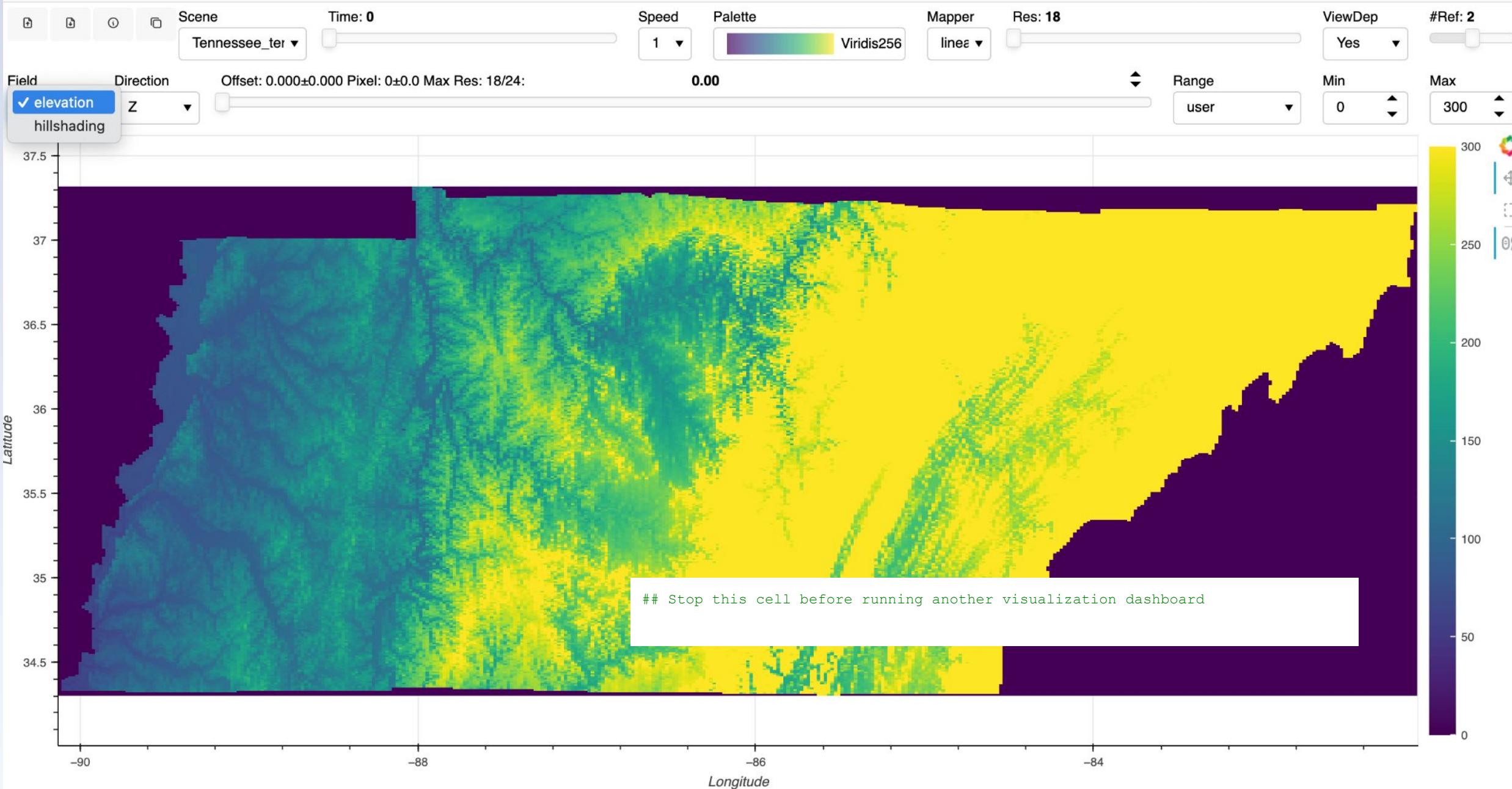
## Option B

From Seal Storage

## Step 4: Interactive Visualization & Analysis

Launch dashboard for interacting with large-scale data to access subregions of the original dataset for ad hoc analysis.





# Step 4: Interactive Visualization & Analysis

Remotely **access** large datasets, **zoom** into specific areas, **select** and **crop** subregions of interest, **save** data locally in a Python-compatible format, and **analyze** the data for scientific discovery.

Step 4 provides **two options** to obtain data and collect the IDX files

**Option A**

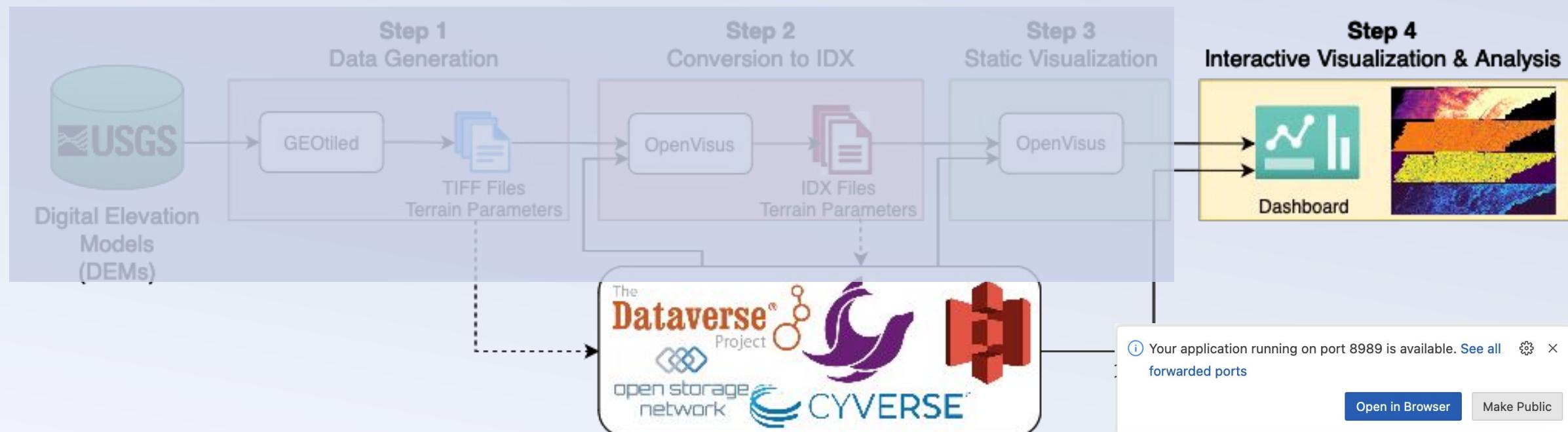
From local storage

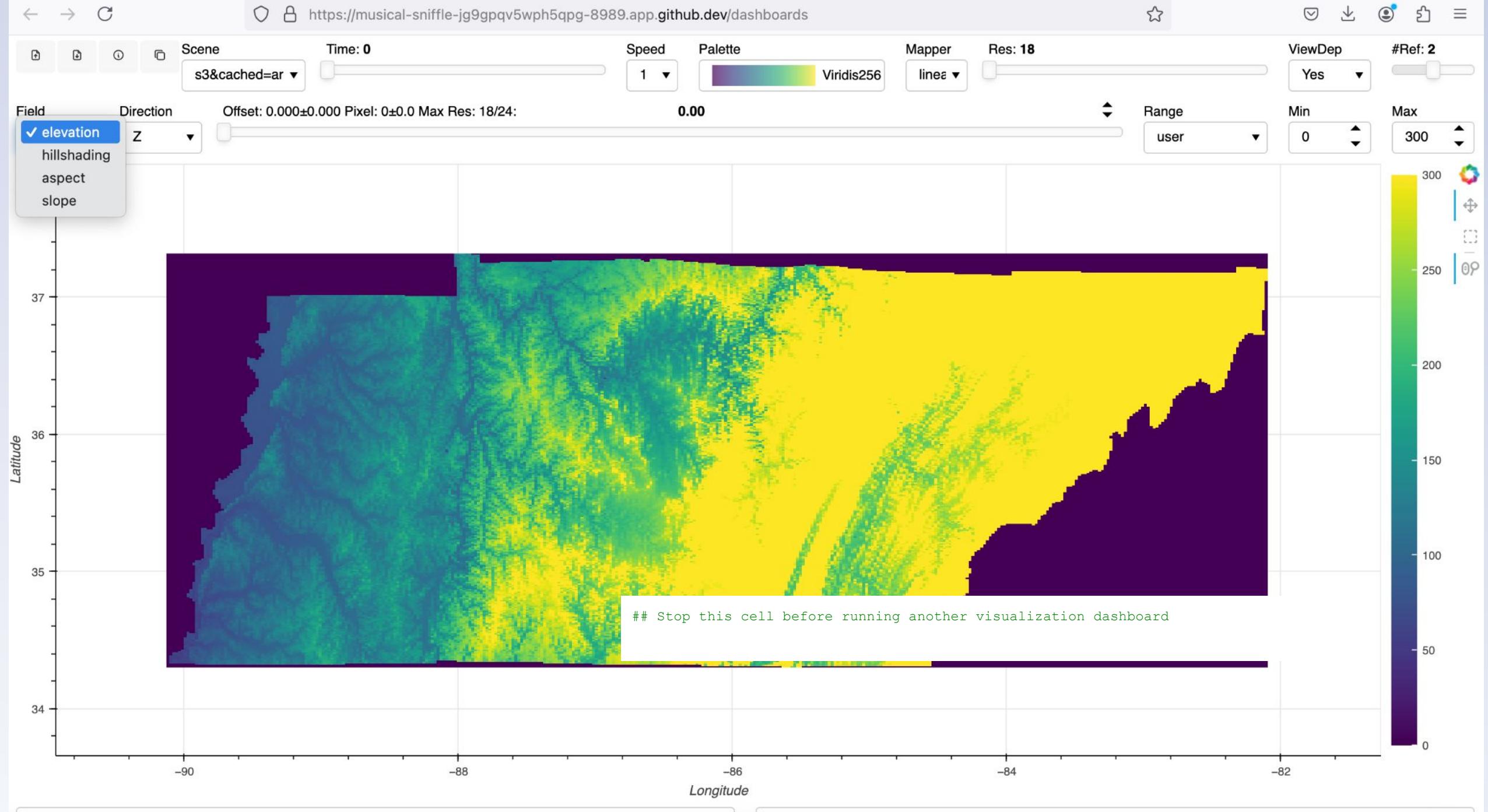
**Option B**

From Seal Storage

## Step 4: Interactive Visualization & Analysis

Launch dashboard for interacting with large-scale data to access subregions of the original dataset for ad hoc analysis.



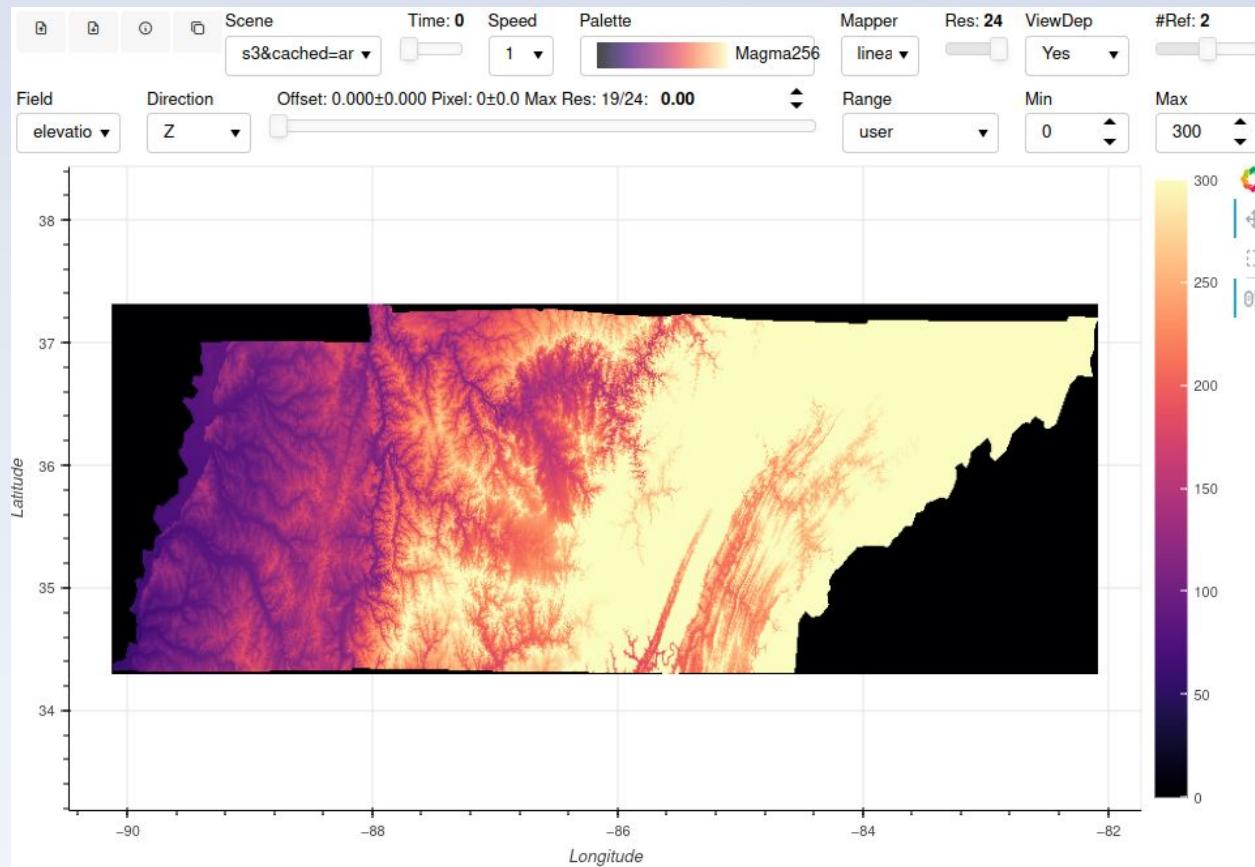




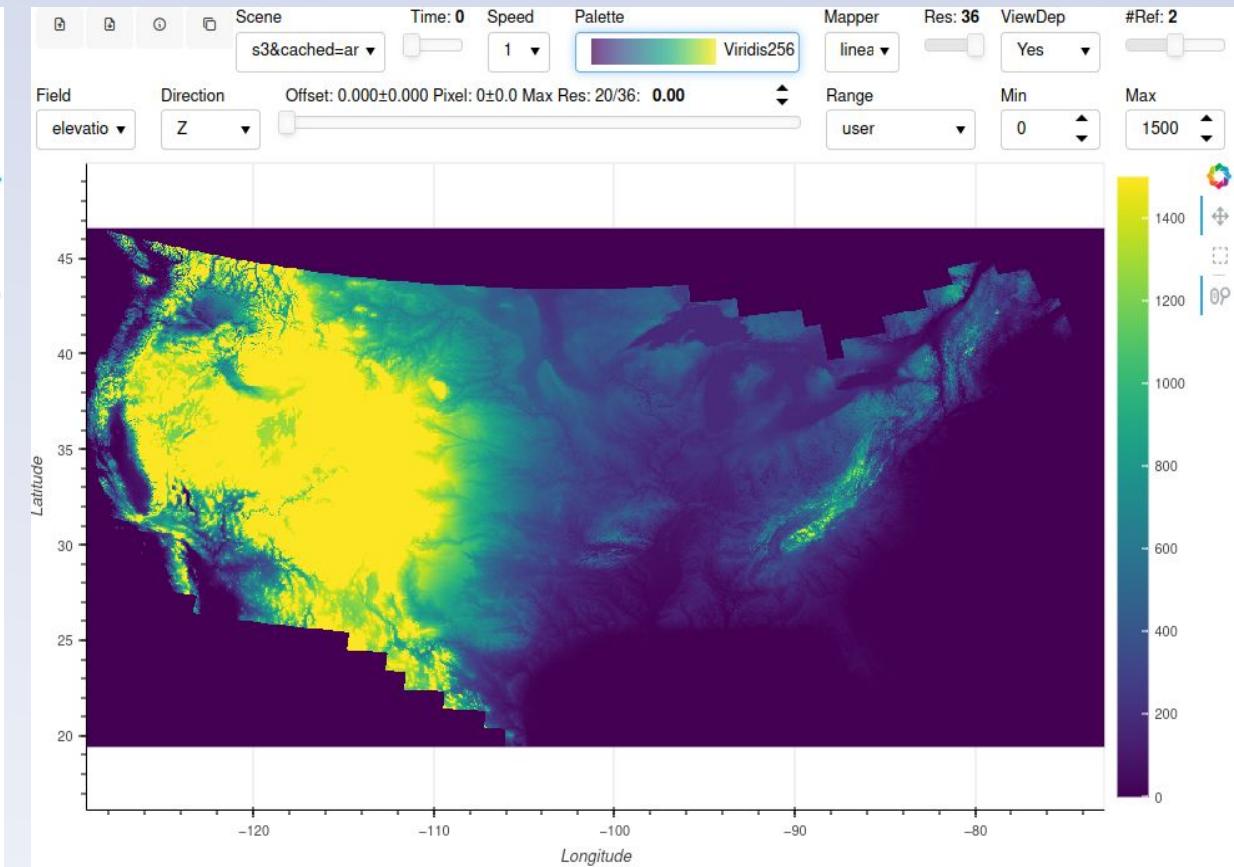
# Step 4: Geographical Regions

Visualize and analyze two geographical regions at 30 m resolution

State of Tennessee - 200 MB

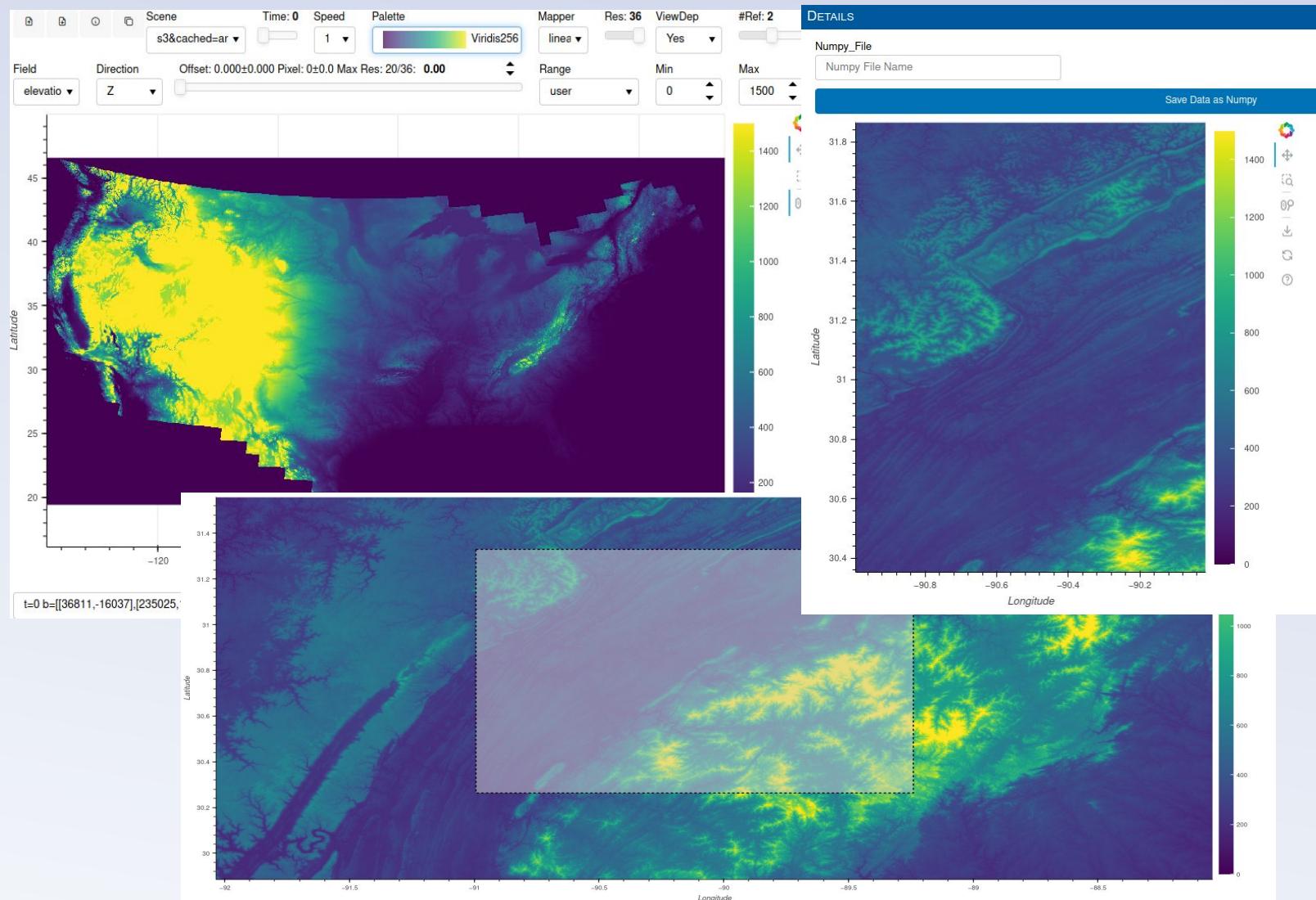


Contiguous United States (CONUS) - 200 GB

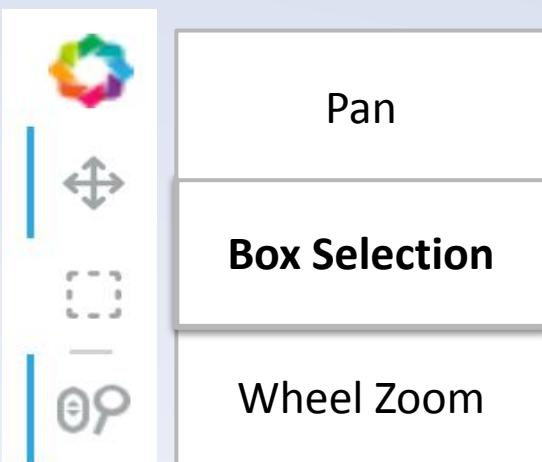




# Step 4: Analysis of Large-Scale Subregions



- Visualize large-scale data **remotely**
- Select and explore subregions
- Save the subregions of interest **locally** in file `zoom-conus-r01`





# Hands on Exercise: Exploring your Subregion Data

- (1) Load the downloaded subregion of interest in your local machine
- (2) Compute min, max, average and std elevation
- (3) Set the color bar to reflect the range of displayed data, from the minimum to the maximum value, providing a more accurate visual representation of the data

Notebook for Exploring Cropped Subregions

After successfully running the tutorial notebook, you can use this jupyter notebook to read and explore the cropped subregion of interest. We present you with two functions to load the data and to statically visualize it. You can expand the analysis of your selected data as required.

Preparing your Environment

The following cell prepares the environment necessary for reading and plotting the data. Upon completion, a message will be displayed to notify you that the cell execution has finished.

```
[1]: import numpy as np
import matplotlib.pyplot as plt
print("You have successfully prepared your environment.")

You have successfully prepared your environment.
```

Enter the name of your Subregion File

Enter the name of the downloaded file.

```
[3]: data_file = "data3.npz"
print("You have successfully named your data file.")

You have successfully named your data file.
```

Reading the Data in the Subregion File

The following cell loads the data and extracts the coordinates and terrain parameter value.

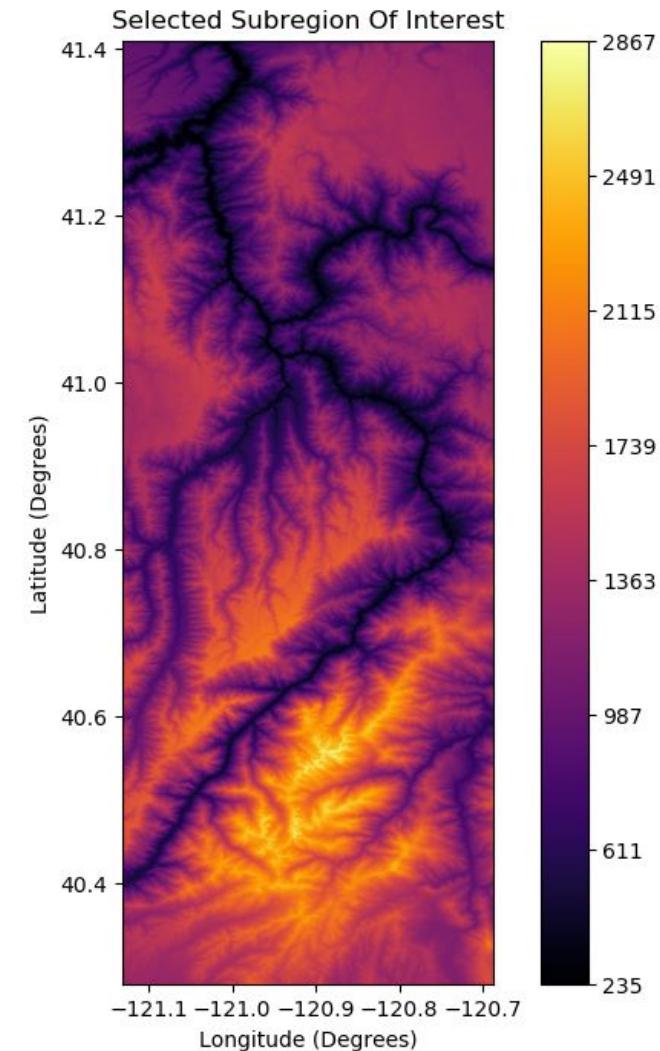
```
[4]: data=np.load(data_file)
data
actual_data=data['data']
metadata=data['lon_lat']
print("You have successfully loaded your data and metadata.")

You have successfully loaded your data and metadata.
```

Visualizing the Subregion Data

The following cell plots the subregion.

```
[5]: cmap_instance = plt.get_cmap("inferno")
lat_min=metadata[0][0]
lat_max=metadata[0][1]
lon_min=metadata[1][0]
lon_max=metadata[1][1]
fig, axes = plt.subplots(1, 1, figsize=(10, 8))
axes.set_xlim(lat_min, lat_max)
axes.set_ylim(lon_min, lon_max)
axes.set_title("Selected Subregion Of Interest")
```





# Hands on Exercise: Exploring your Subregion Data

- (1) Load the downloaded subregion of interest in your local machine
- (2) Compute min, max, average and std elevation
- (3) Set the color bar to reflect the range of displayed data, from the minimum to the maximum value, providing a more accurate visual representation of the data



# Discussion and Open Questions

**Construct a modular workflow** that combines your application components with NSDF services

Can you think of an application that is modular? Can you leverage its APIs?

**Can the application take advantage of the NSDF services?**

**Upload, download, and stream data to and from public and private storage** solutions

How large is the application's data? How do you access, share, and store the data?

**Can the data take advantage of private and public storage?**

**Deploy the NSDF dashboard for large-scale data access, visualization, and analysis**

What type of analysis do you perform on the data?

**Can your research take advantage of an interactive dashboard?**



National Science Data Fabric



SDSC

IBM





# Tutorial Links



[NSDF](#)



[GEOtiled](#)



[SOMOSPIE](#)



[OpenVisus](#)