# The National Data Platform (NDP): Democratizing Data and Responsible Artificial Intelligence

## Manish Parashar

Director, Scientific Computing & Imaging (SCI) Institute

Chair in Computational Science and Engineering

Presidential Professor, Kahlert School of Computing
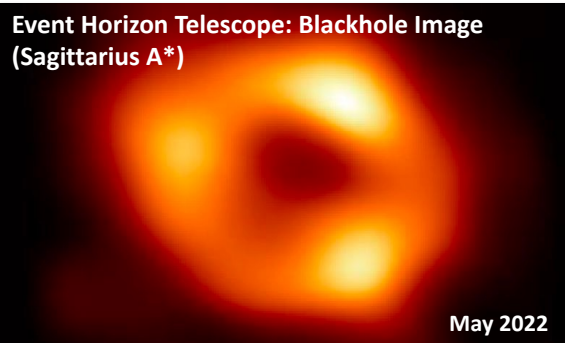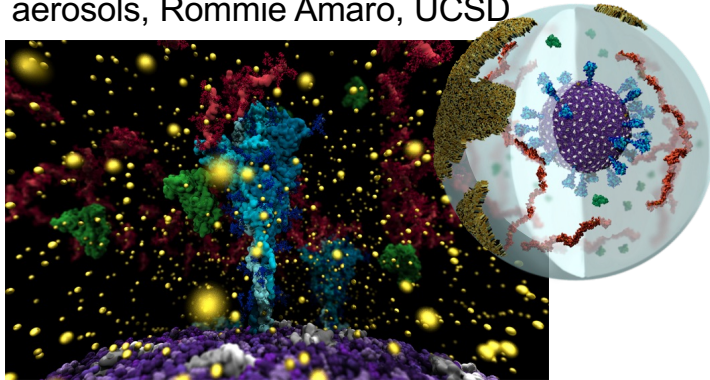
**6NRP**
San Diego, CA
January 30, 2025

**Scientific Computing and Imaging (SCI) Institute**          **One-U Responsible AI Initiative**

# Science / Society Transformed by Data, Cyberinfrastructure

Modeling of the delta virus inside respiratory aerosols, Rommie Amaro, UCSD



Cyberinfrastructure is a key enabler of discoveries & innovations

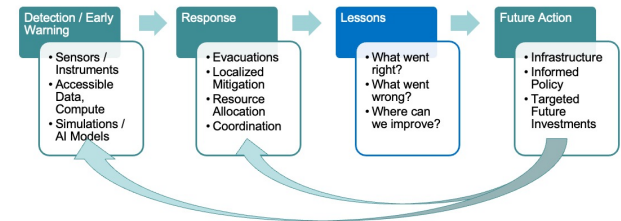Data-driven Urgent Science



Earthquakes & Tsunamis
Extreme Weather
Cyber-Attacks
Pandemics
Industrial Disasters

Event Horizon Telescope: Blackhole Image (Sagittarius A*)

May 2022



Detection / Early Warning
- Sensors / Instruments
- Accessible Data, Compute
- Simulations / AI Models

Response
- Evacuations
- Localized Mitigation
- Resource Allocation
- Coordination

Lessons
- What went right?
- What went wrong?
- Where can we improve?

Future Action
- Infrastructure
- Informed Policy
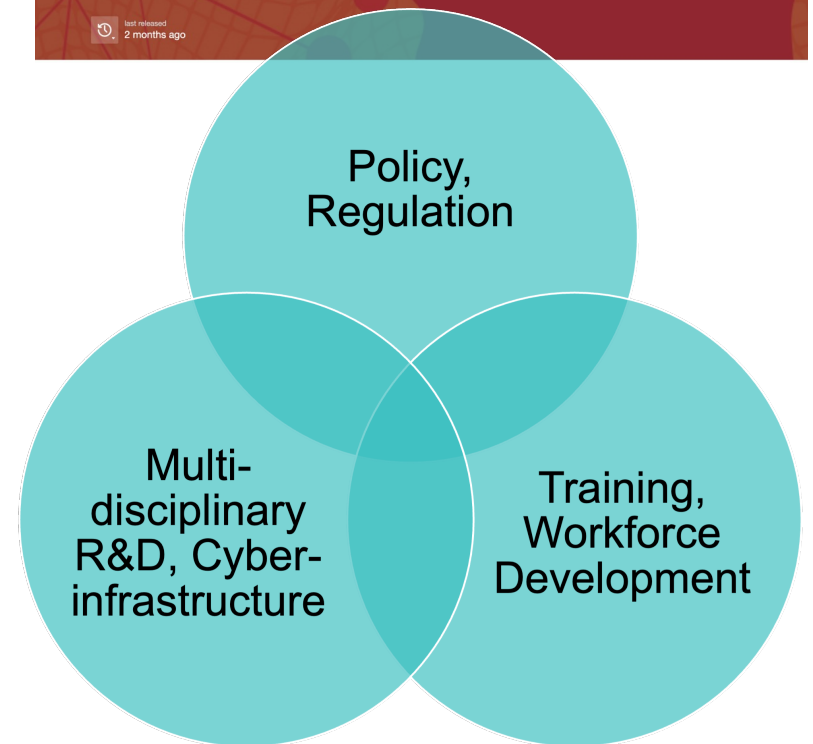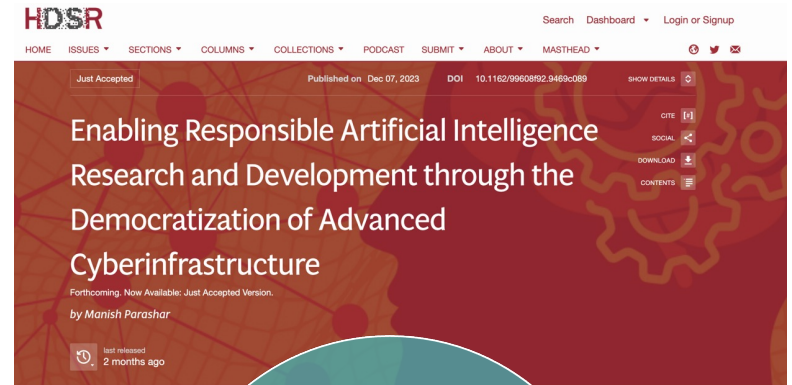- Targeted Future Investments

# The Transcendence of AI

# Democratizing Responsible Data, AI is important

- Key characteristics for responsible AI (NIST) *validity, safety security, accountability, privacy enhancement, fairness, and explainability.*

- The quality and impact of research and the pace of innovation are linked to the diversity of the contributions.

- Especially true for AI-enabled research
  - In case of AI, quality depends on who is developing and use AI, and where the data coming from
  - Greater inclusivity in contribution to research and development increases the diversity of approaches and the fairness of the results.

- Many barriers: awareness, ability, access, association, …

# NAIRR: Democratizing the AI R&D Ecosystem

**Goals:** Strengthen and democratize the U.S. AI Innovation ecosystem in a way that protects privacy, civil rights, and civil liberties.
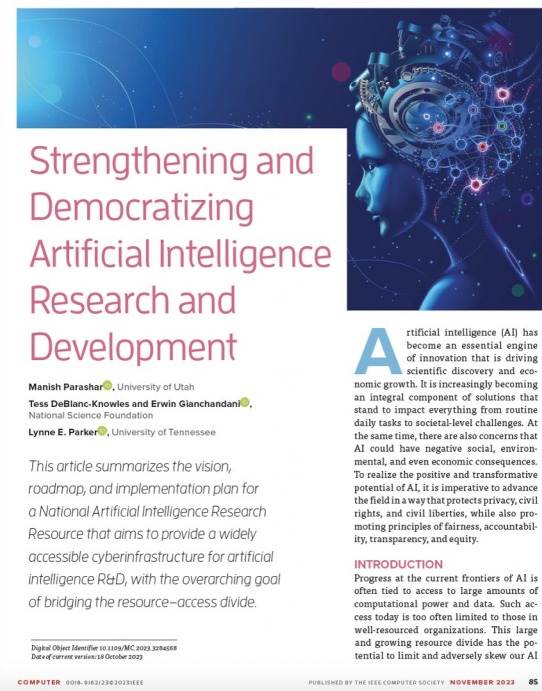
Spur **innovation**

Increase the **diversity** of talent in AI

Improve U.S. **capacity** for AI R&D

Advance **trustworthy AI**



*Computer*, vol. 56, no. 11, pp. 85-90, Nov. 2023, doi: 10.1109/MC.2023.3284568.



https://www.ai.gov/nairrtf/

# National Data Platform - NDP
## Services for Equitable Open Access to Data

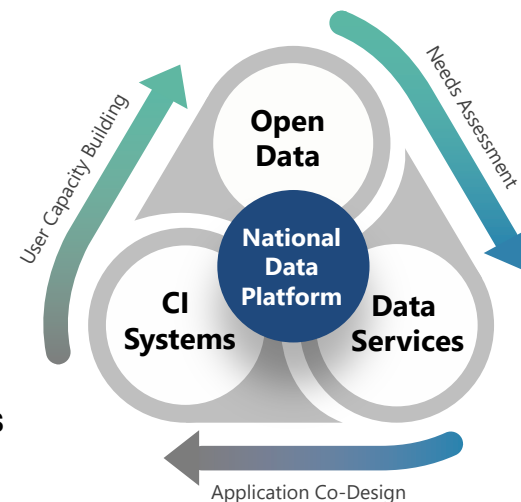**NATIONAL DATA PLATFORM**
Bridging the Data Gaps for AI

A **federated** and **extensible** data ecosystem to promote innovation and collaboration through the equitable use of data leveraging existing and future national cyberinfrastructure capabilities.

https://www.nationaldataplatform.org/

## FOCUS AREAS:

- **Platform** for data-enabled and AI-integrated workflows
  - Facilitates data registration and discovery via a **centralized hub**
  - Democratizes data access and use via **distributed points of presence**
  - Cultivates resources for **classroom education** and **data challenges**
  - Assists research and learning through **personalized workspaces**

- **Applications** in climate and AI with data diverse scientific data repositories including NSF facilities, NAIRR, NASA, USGS, NOAA and USDA

- **Partnerships** to foster scientific discovery, decision-making, policy formation and societal impact

University of Colorado Boulder



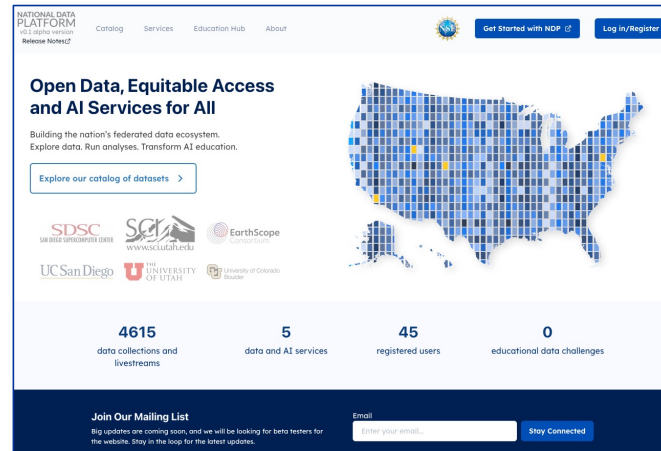UC San Diego · THE UNIVERSITY OF UTAH · University of Colorado Boulder · SAN DIEGO SUPERCOMPUTER CENTER · SCI www.sci.utah.edu · EarthScope Consortium

**Centralized portal** for discovery through collaborative workspaces for research and education
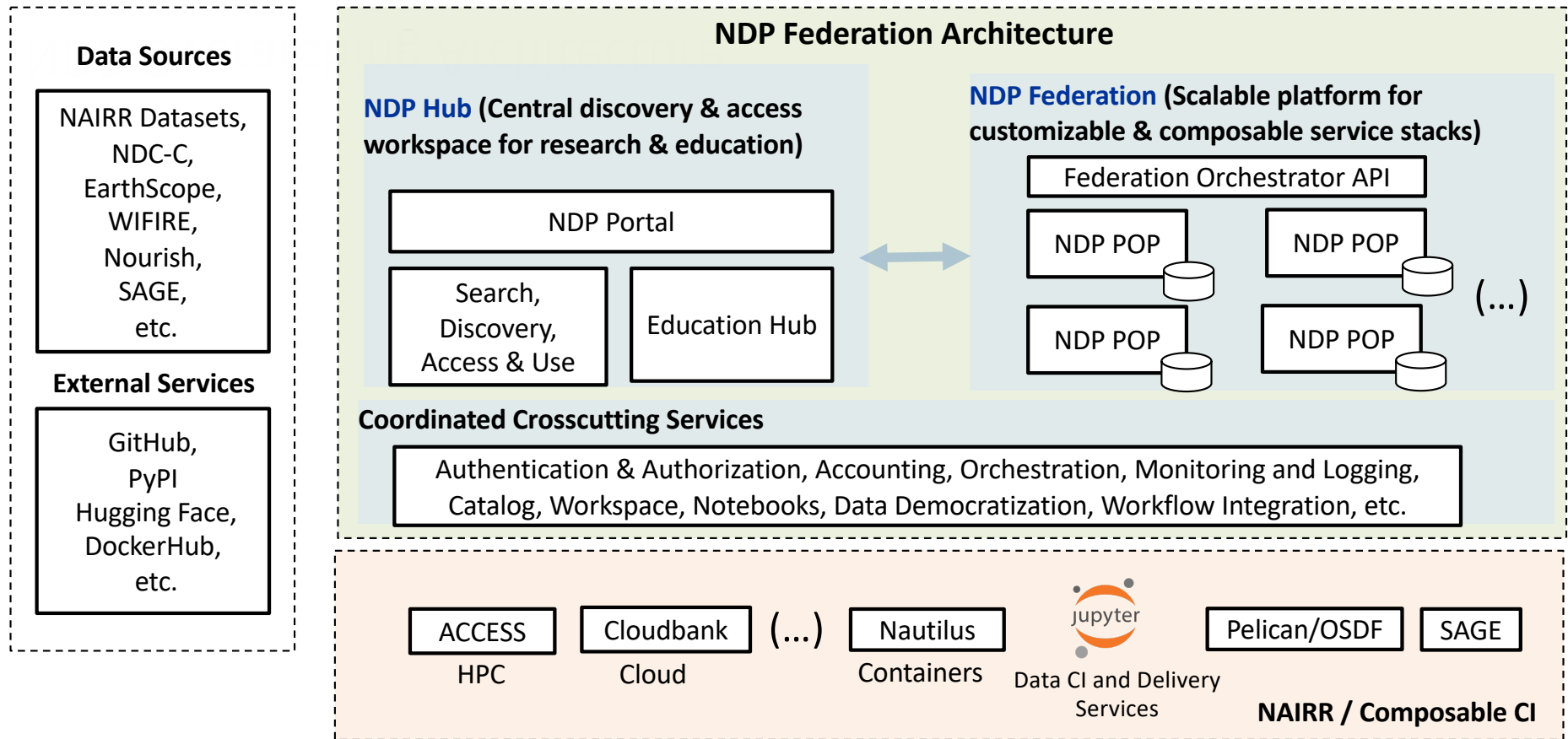
**NDP Federation**

A scalable **platform** for developing and deploying services at **distributed points of presence**
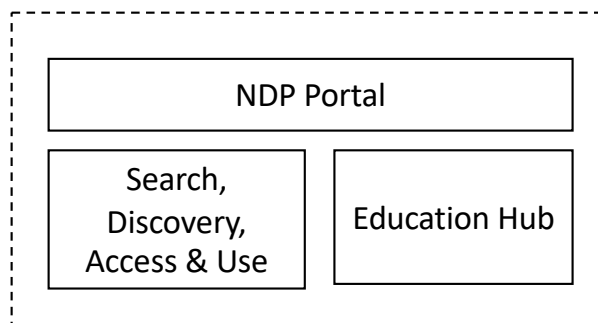
NATIONAL DATA PLATFORM

# NDP Overarching Architecture

## Data Sources

NAIRR Datasets,
NDC-C,
EarthScope,
WIFIRE,
Nourish,
SAGE,
etc.

## External Services

GitHub,
PyPI
Hugging Face,
DockerHub,
etc.

## NDP Federation Architecture

### NDP Hub (Central discovery & access workspace for research & education)

NDP Portal

Search, Discovery, Access & Use

Education Hub

### NDP Federation (Scalable platform for customizable & composable service stacks)

Federation Orchestrator API

NDP POP

NDP POP

(...)

NDP POP

NDP POP

### Coordinated Crosscutting Services

Authentication & Authorization, Accounting, Orchestration, Monitoring and Logging, Catalog, Workspace, Notebooks, Data Democratization, Workflow Integration, etc.

ACCESS
HPC

Cloudbank
Cloud

(...)

Nautilus
Containers

jupyter
Data CI and Delivery Services

Pelican/OSDF

SAGE

NAIRR / Composable CI

# NDP Hub: Central discovery & access workspace for research & education

**NDP Hub**

| NDP Portal |
|---|

| Search, Discovery, Access & Use | Education Hub |
|---|---|

- NDP Portal (point of access)

  **https://nationaldataplatform.org**

- Metadata registration and indexing
  - Contributing organizations
  - Harvested metadata from NDP POPs
- Data search
  - String and conceptual search
  - Open Knowledge graphs / via LLMs

**NDP Standard Services**

Public:
- Extensible Data Catalog and Search Services
- Education Hub Informal Learning Modules

Login-enabled:
- Keycloak Role-Based Access Service
- User Workspaces
- AI Gateway with Custom JupyterHub Service
- Data Catalog and OKN Ingestion
- External Model Ingestion
- Data Exploration Services
- MLFlow Dashboard Service
- Education Hub Classroom
- Education Hub Challenge
- Democratizing Data Dashboard

**Hub Capabilities Under Development**

- Sage Data and Edge Code Integration Service
- Service Catalog and Discovery Service
- Educational Hub Expansion
- Streaming Data Services
- Pelican Registration Service
- Integrated Workflows

**Planned Future Work**

- OKN Integration
- Data Curation
- Data Subsetting
- Data Provenance
- Educational Toolkits
- Open Science Chain Provenance Service
- Gateway Services

# NDP Hub: Data Search and Discovery



**Current Capabilities:**
- Search capabilities to include not just text in metadata and ontology concepts but also time and location data.
- Ability to search time and time ranges within the data, such as from "27 September 2020" to "24 January 2021."
- Location-based searches can now be combined using specific location names (e.g., "San Luis Obispo") or boundary polygons.
- Support free-text search across "all metadata" without specifying particular fields.
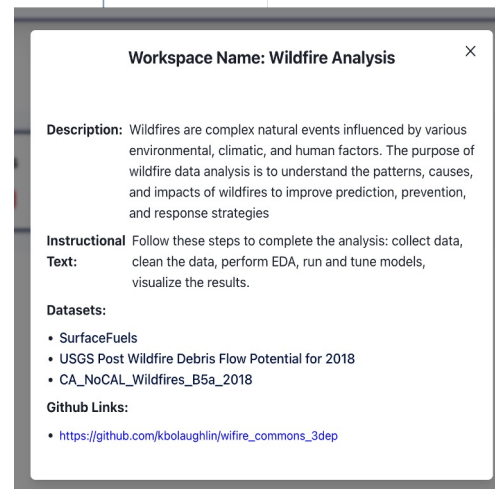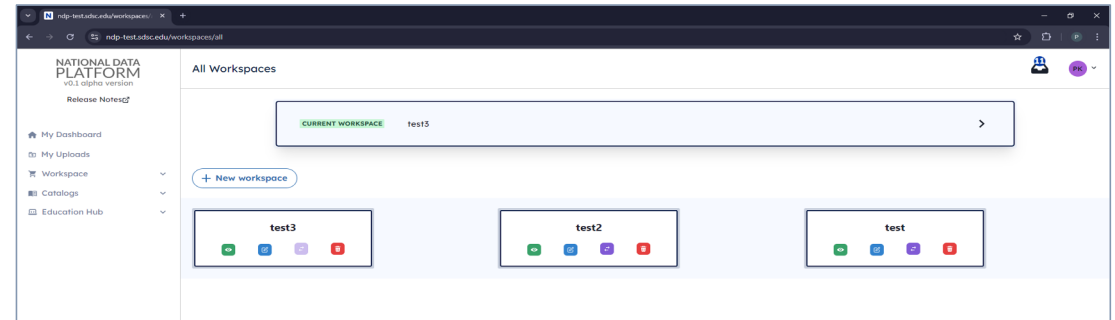- Utilize Lucene, a popular search syntax, to improve search functionality.

**Ongoing Work :**
- Extract entity annotations from the metadata text and integrate them with the ontology to enhance search functionality.
- Create a vector store and develop a search pipeline that handles queries in natural language.
- Optimize the system's performance to ensure fast and accurate retrieval of relevant information.

# NDP Workspaces (Version 1 – September 2024)

**Goal :** Craft persistent and customizable workspaces with datasets and services to launch into a sandbox

- Create customized workspaces for varied use cases
- Search and add datasets to use in sandbox (HPC Env)
- Add github links for file access
- Launch packaged workspace into sandbox



**NATIONAL DATA PLATFORM**
v0.1 alpha version
Release Notes

- My Dashboard
- My Uploads
- Workspace
- Catalogs
- Education Hub

All Workspaces

CURRENT WORKSPACE    test3

+ New workspace

test3    test2    test

**Workspace Name: Wildfire Analysis** ✕

**Description:** Wildfires are complex natural events influenced by various environmental, climatic, and human factors. The purpose of wildfire data analysis is to understand the patterns, causes, and impacts of wildfires to improve prediction, prevention, and response strategies

**Instructional Text:** Follow these steps to complete the analysis: collect data, clean the data, perform EDA, run and tune models, visualize the results.

**Datasets:**
- SurfaceFuels
- USGS Post Wildfire Debris Flow Potential for 2018
- CA_NoCAL_Wildfires_B5a_2018

**Github Links:**
- https://github.com/kbolaughlin/wifire_commons_3dep

Users can:
- view all their workspaces
- create new workspaces by clicking on the "New Workspace" button
- use workspace action buttons to preview, edit, switch and delete
- add datasets to their current workspace from the catalog page.

SDSC SAN DIEGO SUPERCOMPUTER CENTER — SCI www.sci.utah.edu — UTAH U — University of Colorado Boulder — EarthScope Consortium — http://www.nationaldataplatform.org — UC San Diego HALICIOĞLU DATA SCIENCE INSTITUTE

# NDP JupyterHub (Sandbox)

A compute environment for data analysis, machine learning training or any other computational tasks, built on top of NRP (Nautilus) cluster. Different datasets and tasks will require powerful compute resources (CPUs, GPUs, memory), which user can select and use seamlessly.



✓ Integrated with NDP Single-Sign On

✓ Select your compute resources from NRP pool

✓ Select previously created image (environment) or bring yours

- Integrated with File Manager extension
- Loads data from your workspaces (datasets and github resources)
- Change your workspaces content and refresh in JupyterHub to get updates
- Download all or selected resources into your storage for further analysis

# NDP Data POP: Distributed Points of Presence with Customizable, Composable Service Stacks
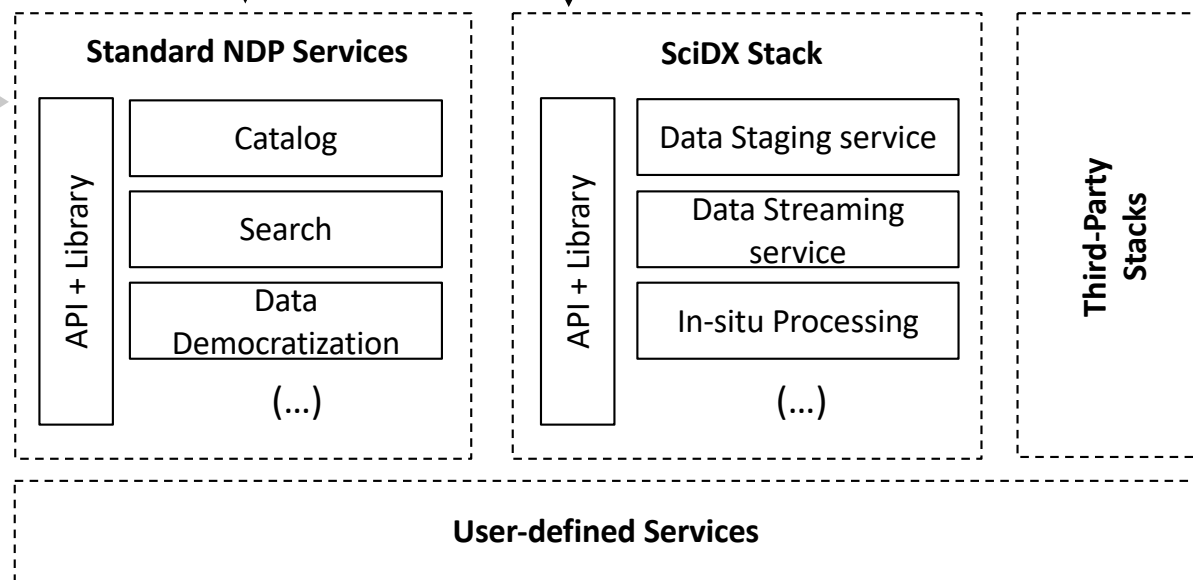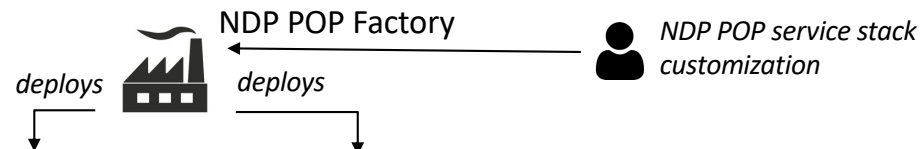
# Science Data Exchanges (SciDx) Services
*A customizable Data-Pop software stack for in-situ data access & processing*
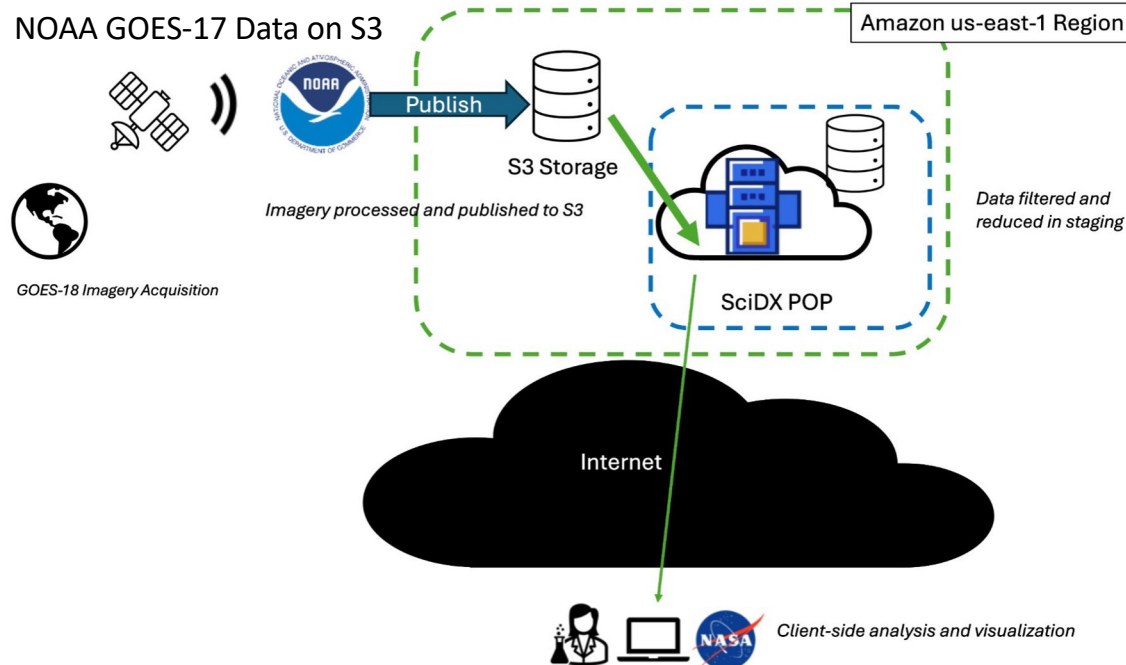
## SciDx Staging Services

- Transient resources for in-situ (close to the data) data processing and access
  - High-performance in-memory processing
  - Server-side data transformations (e.g., sub-setting, reduction, user-defined analysis, etc.)
  - Caching/sharing of data, results, and data-products
  - Registration of data-triggers
- Efficient management of data in-motion
  - Streamline workflows; minimize data transfers
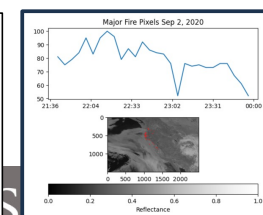  - Perform ETL operations at data source

# SciDx Staging Service: Wildfire Monitoring Usecase

- Monitor fire hotspots based on satellite data that updates every 5 minutes
- Not interested in the entire data product, just pixels that reach severity threshold
- Per-pixel evaluation as a user-defined transformation is performed on each new data update
- The user subscribes to the results of the transformation
- Reduction in data cost, latency, time to solution

NOAA GOES-17 Data on S3

Amazon us-east-1 Region

Publish

S3 Storage

Imagery processed and published to S3

GOES-18 Imagery Acquisition

Data filtered and reduced in staging

SciDX POP

Internet

Client-side analysis and visualization

```
result =
client.query_array(source='goes18-radc',
            var_name='Rad',
            lb=(0,2500),
            ub=(2499,4999),
            timestamp='2024-08-02T00:35:00',
            time_direction=PAST)
```

Major Fire Pixels Sep 2, 2020

# Science Data Exchanges (SciDx) Services

*A customizable Data-Pop software stack for in-situ data access & processing*
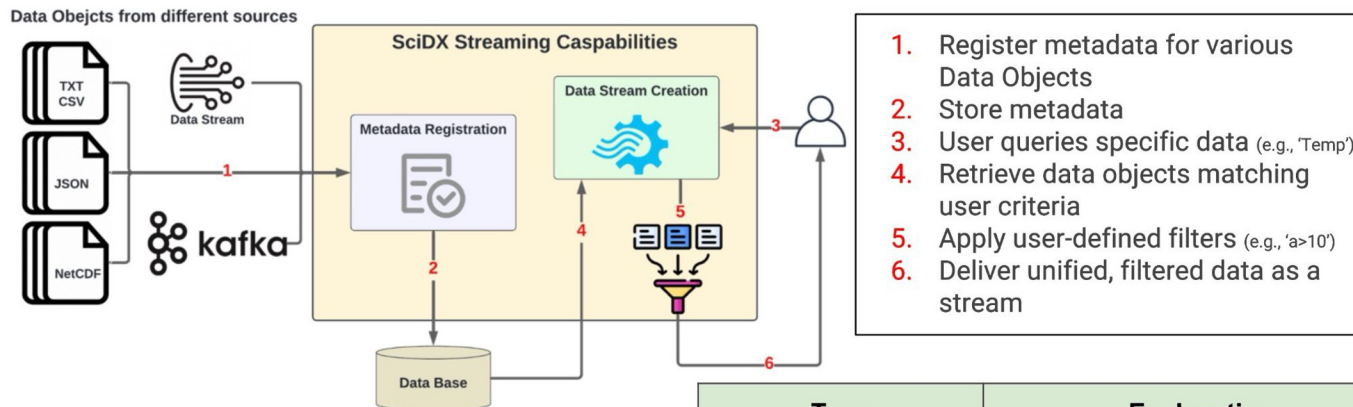
## SciDx Staging Services

- Transient resources for in-situ (close to the data) data processing and access
  - High-performance in-memory processing
  - Server-side data transformations (e.g., sub-setting, reduction, user-defined analysis, etc.)
  - Caching/sharing of data, results, and data-products
  - Registration of data-triggers
- Efficient management of data in-motion
  - Streamline workflows; minimize data transfers
  - Perform ETL operations at data source

## SciDx Streaming Service

- Streams registration, curation/archival for discovery and access
- User-defined operations/filters on streaming; containerized execution
- Combine streaming data with archived/playback data
- Mechanism for online data product generation (i.e., new data streams)
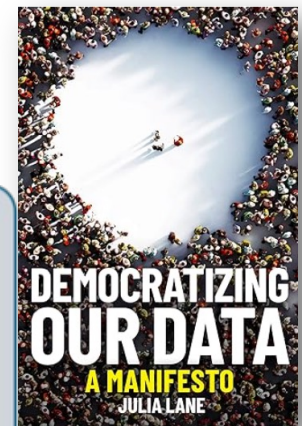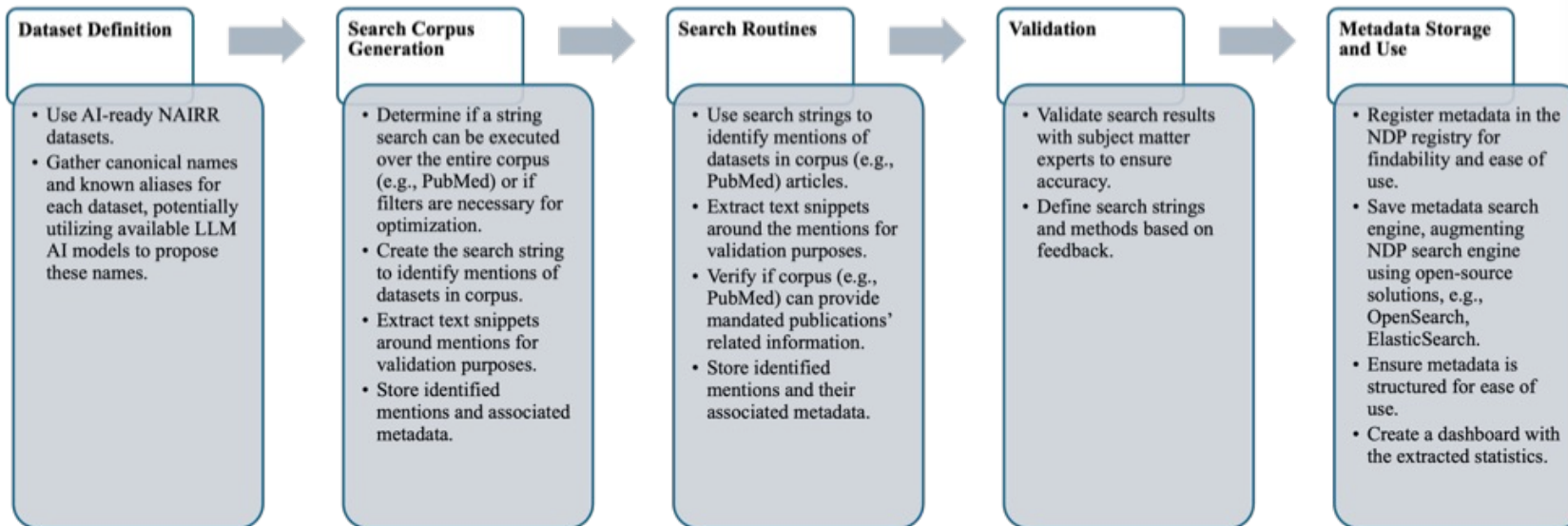
# SciDx Streaming Service



1. Register metadata for various Data Objects
2. Store metadata
3. User queries specific data (e.g., 'Temp')
4. Retrieve data objects matching user criteria
5. Apply user-defined filters (e.g., 'a>10')
6. Deliver unified, filtered data as a stream

| Type | Explanation | Example |
|---|---|---|
| Column Comparisons | Column-to-column comparisons | `x > y` |
| Mathematical Operations | Addition, subtraction, multiplication and division | `x > 10*y` |
| IN Operator | Check if values are in a list | `station IN ['A', 'B']` |
| Conditional Logic (IF-THEN-ELSE) | Apply rules based on conditional statements | `IF x > 20 THEN alert = High ELSE y = 10` |
| Logical Operators (AND, OR) | Combine multiple conditions using AND and OR operators | `IF x > 10 OR z = 20 THEN alert = High ELSE alert = Low` |
| Window-Based Filtering | Calculate aggregates (mean, sum, max, min) over sliding windows | `IF window_filter(9, sum, x > 20) THEN alert = High` |

# SciDx: Advanced Search & Discovery

https://democratizingdata.ai/



**Dataset Definition**
- Use AI-ready NAIRR datasets.
- Gather canonical names and known aliases for each dataset, potentially utilizing available LLM AI models to propose these names.

**Search Corpus Generation**
- Determine if a string search can be executed over the entire corpus (e.g., PubMed) or if filters are necessary for optimization.
- Create the search string to identify mentions of datasets in corpus.
- Extract text snippets around mentions for validation purposes.
- Store identified mentions and associated metadata.

**Search Routines**
- Use search strings to identify mentions of datasets in corpus (e.g., PubMed) articles.
- Extract text snippets around the mentions for validation purposes.
- Verify if corpus (e.g., PubMed) can provide mandated publications' related information.
- Store identified mentions and their associated metadata.

**Validation**
- Validate search results with subject matter experts to ensure accuracy.
- Define search strings and methods based on feedback.

**Metadata Storage and Use**
- Register metadata in the NDP registry for findability and ease of use.
- Save metadata search engine, augmenting NDP search engine using open-source solutions, e.g., OpenSearch, ElasticSearch.
- Ensure metadata is structured for ease of use.
- Create a dashboard with the extracted statistics.

DEMOCRATIZING OUR DATA
A MANIFESTO
JULIA LANE

# NDP+NRP: Use the NDP widget to import datasets and conduct analysis within NRP.



Students select a module to work on. They load it to JHub, cloning the attached repository and loading the data using the NDP Widget

A reminder for students to save their work in a persistent storage directory

# Example NDP-NAIRR AI in Science Workflow

Data Acquisition → Registration, Indexing, Discovery → Workspace Definition → Data-driven, AI/ML-based workflows → Product Generation, Curation, Sharing, Archival

- Data and Models are identified as part of the Open NAIRR Resources.
- Resources are collected from HuggingFace

- Data and Models are registered into NDP catalog (CKAN)
- Data origin is created in OSDF to optimize data transfer

- Data and Models are included into user's workspace, along with the necessary libraries, services and files to work on a new project.

- Analysis and AI/ML workflow is supported by AI Gateway (JupyterHub), using NRP's Nautilus.
- High Performance processing for new resource(s) development (Models, Data).

- Final products pushed to OSDF/HuggingFace/GitHub and registered into NDP's catalog .
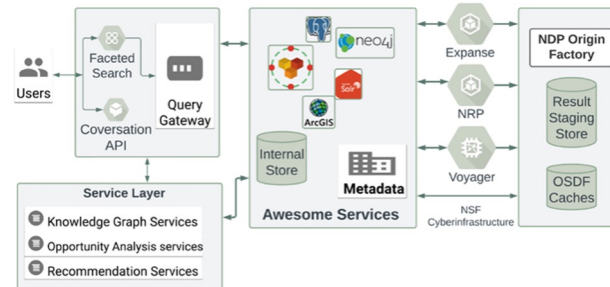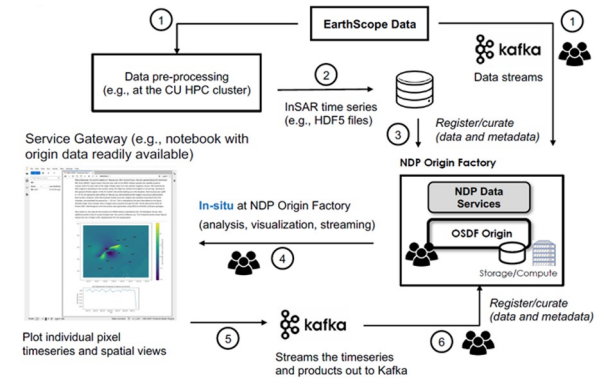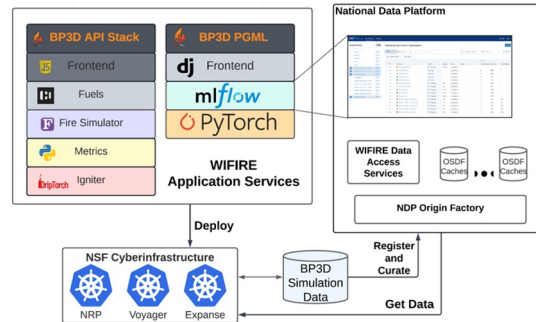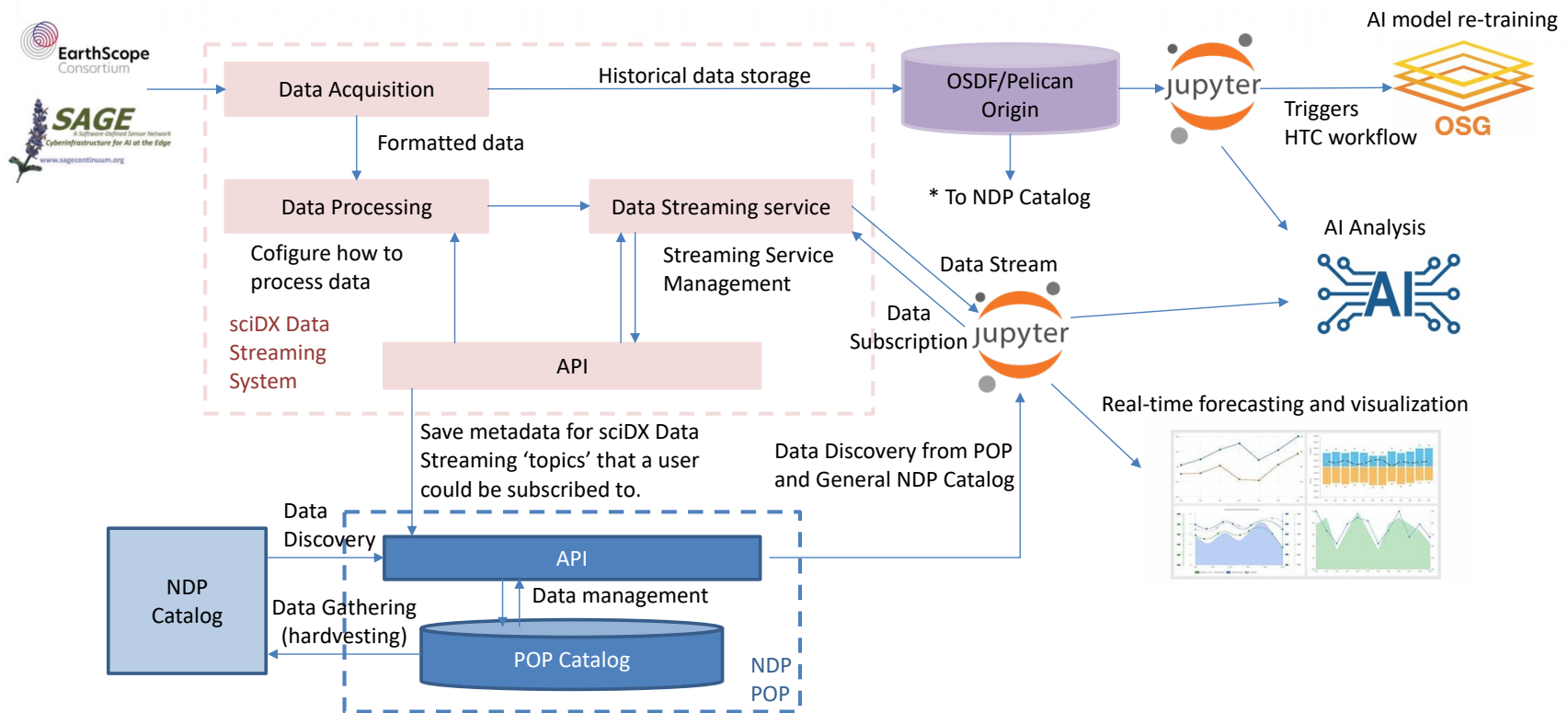


http://www.nationaldataplatform.org

# Case Studies for Generalizable Workflows

- **Representative examples** of important patterns that exist in science today for working with
  - large datasets
  - streaming data from facilities
  - graph data from open knowledge networks

- Implemented as production-quality specialized value-added services

- Domains of wildland fire, earthquakes, and food security

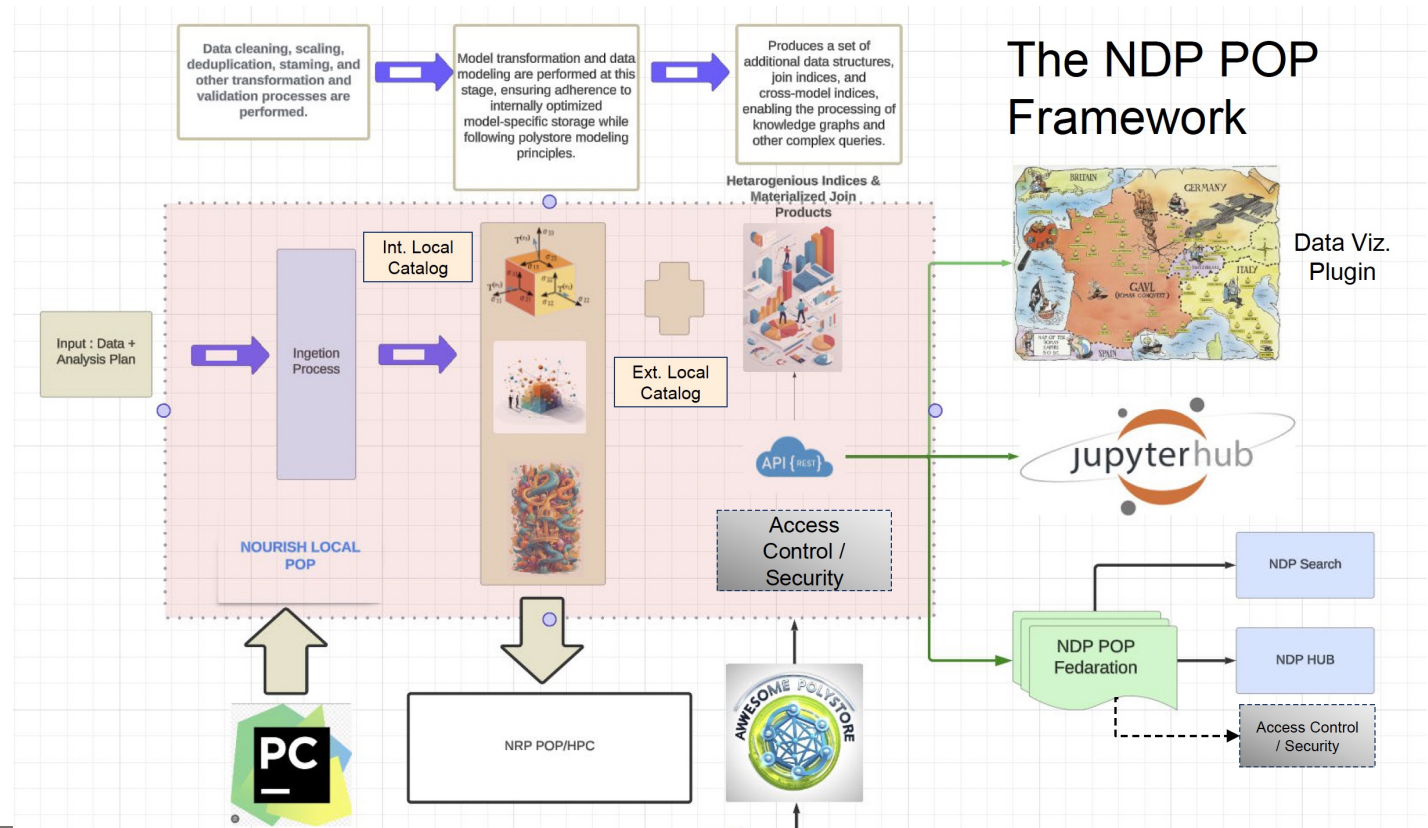- Will be generalized for replication by external communities.

# EarthScope/SAGE data streaming/analysis enabled by NDP POP

- Real-time high-precision GNSS stations and SAGE data streams

# NDP + NDR: Nourish NDF POP



The NDP POP Framework

Data Viz. Plugin

Extends SDSC's AWESOME platform

**NATIONAL DATA PLATFORM**

**Bridging the Data Gaps for AI**

UC San Diego

UTAH U

University of Colorado Boulder

SDSC SAN DIEGO SUPERCOMPUTER CENTER

SCI www.sci.utah.edu

EarthScope Consortium

http://www.nationaldataplatform.org

SDSC SAN DIEGO SUPERCOMPUTER CENTER

İlkay Altıntaş, PhD (ialtintas@ucsd.edu )

UC San Diego
HALICIOĞLU DATA SCIENCE INSTITUTE
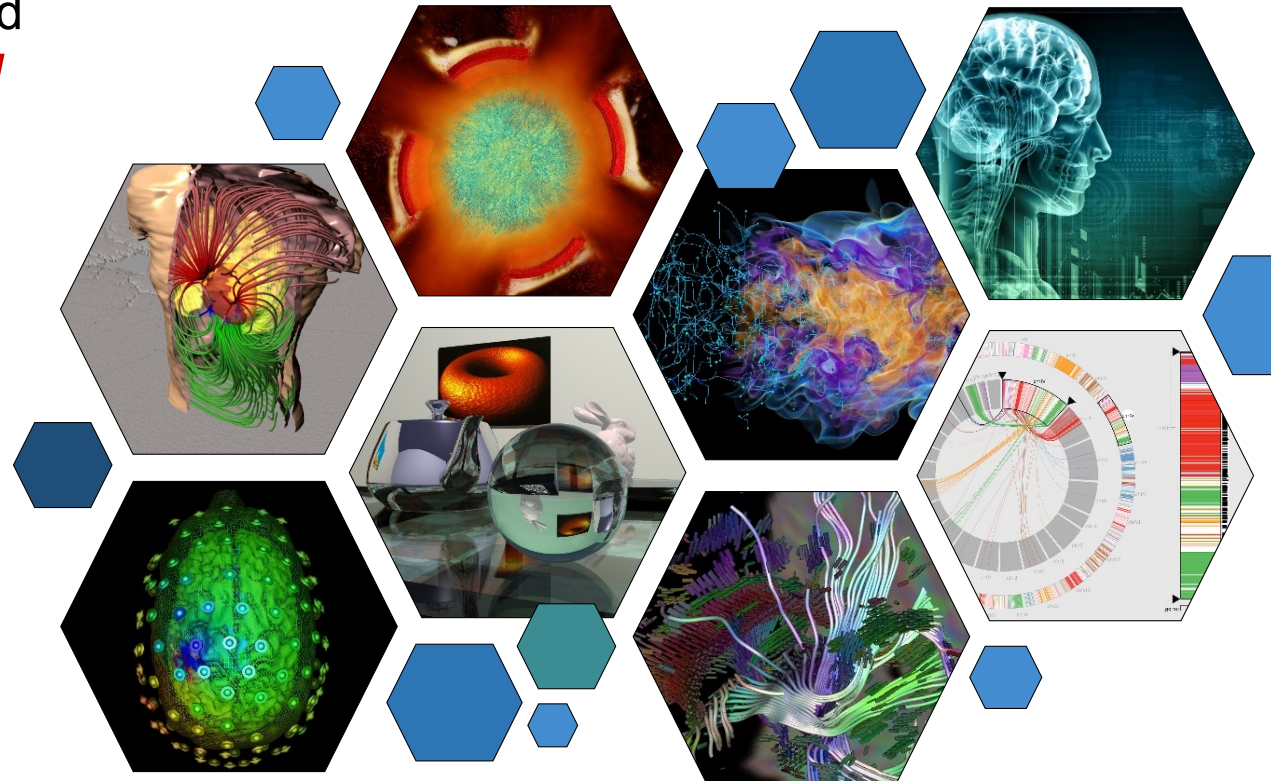
University of Colorado Boulder

# SCI Institute

Transformation of science and society through ***translational research and innovation***

- Inter/transdisciplinary, collaborative, convergent

- Core strengths in: Visualization & imaging; Scalable analytics; Advanced computing & data

- Software/system development and distribution integral to our research processes
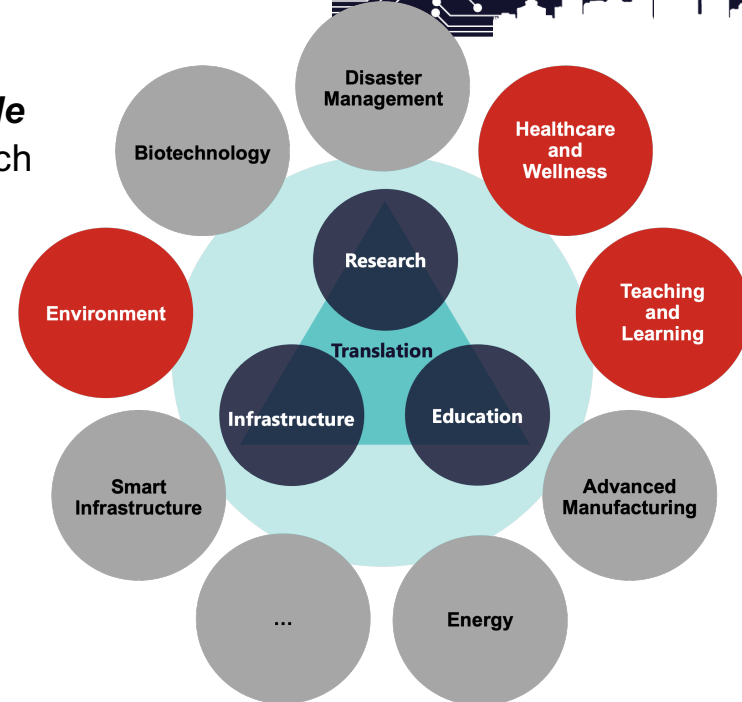
# One-U Responsible AI Initiative

## University of Utah's $100M AI Research Initiative Led by SCI

Responsibly advance *translational AI* to achieve societal good

> Catalyze *transdisciplinary excellence in responsible AI* to bring together the use-inspired/applied AI research and technological expertise, advanced cyberinfrastructure, and translational workforce

Initial research focus on regionally important applications

1. Environment
2. Healthcare, societal wellness, and public services
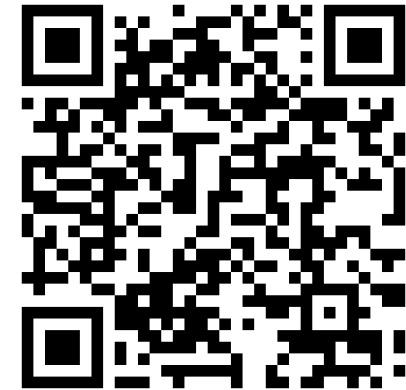3. Future of teaching and learning

# Thank you!

## One-U RAI Opportunities

- **Distinguished Visitors Program:** Supports visits from a few days to a full year for faculty, up to two years for postdoctoral fellows.
- **Postdoctoral Fellows Program:** Supports postdocs in areas related to responsible AI for up to 2 years.

Manish Parashar

Email: manish.parashar@utah.edu

WWW: manishparashar.org / sci.utah.edu / rai.utah.edu

in one-u-responsible-ai-initiative

X @OneU_RAI