

Keynote Talk for the 5th Global Research Platform (GRP) Workshop September 16, 2024 – Osaka, Japan



University of California, San Diego

Chief Data Science Officer & Division Director of Cyberinfrastructure and Convergence Research and Education, San Diego Supercomputer Center

Founding Fellow, Halicioğlu Data Science Institute

Founding Director, Workflows for Data Science Center of Excellence

Founding Director, WIFIRE Lab

Joint Faculty Appointee, Los Alamos National Laboratory



UC San Diego ...
HALICIOĞLU DATA SCIENCE INSTITUTE

SCHOOL OF COMPUTING, INFORMATION AND DATA SCIENCES

UC San Diego

SAN DIEGO SUPERCOMPUT at the UNIVERSITY OF CALIFORNIA S ABOUT SDSC SERVICES SUPPORT RESEARCH & DEVELOPMENT ED Materials Science Researchers Double Up on SDSC, PSC Supercomputers to Discover No **Details about TMDs** Supercomputer simulations provide a better understanding o two-dimensional layered materials showing promise for a var applications - from flexible electronics and spintronics to opto and memory devices. READ MORE Innovate, FOR UC/UCSD Researchers National HPC Users

School of Computing, Information and Data Sciences https://scids.ucsd.edu/

https://www.



UC Regents Approve New School of Computing, Information and Data Sciences at UC San Diego

New school meets critical demand to advance data science and AI innovations and educate workforce of the future



Pioneering Data Science for a

Data-Driven Future

w Does ChatGPT Work? - Event

JULY 18, 2023 - KALEIGH O'MERRY

OCT 8:00 am - 5:00 pm 🗇 12 Swarup Swaminathan, MD | University of Miami Miller

OCT 2:00 pm - 3:00 pm 🗇

Tweets from @HDSIUCSD



Nothing to see here vet

When they Tweet, their Tweets will show up here

cience.ucsd.edu/



UC San Diego HALICIOĞLU DATA SCIENCE INSTITUTE

Cyberinfrastructure and Convergence Research Division @SDSC

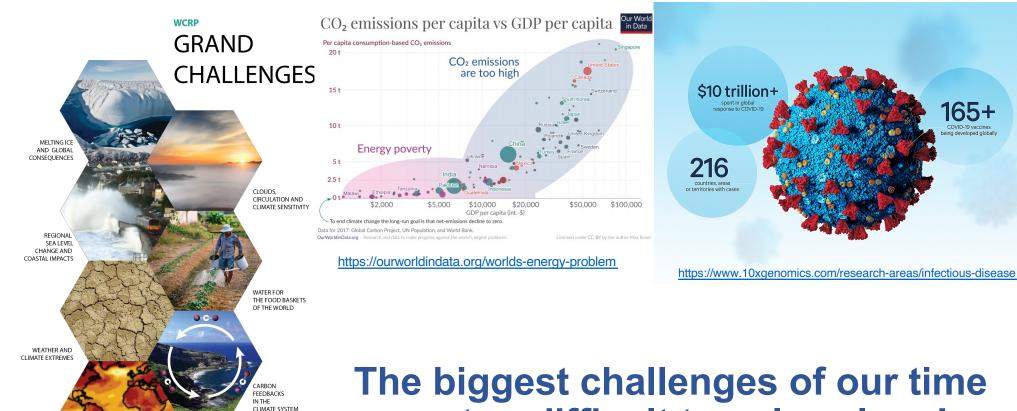
Translating cyberinfrastructure research for impact at scale

- "Big" Data
- Computational Science
- Data Science
- Cyberinfrastructure
- Collaborative Problem Solving
- Convergence Research
- Experiential Education









The biggest challenges of our time are too difficult to solve alone!

https://www.wcrp-climate.org/learn-grand-challenges



CLIMATE



Convergence research is:

driven by a specific and compelling societal problem

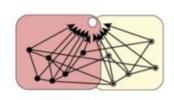
and

works towards integrating innovative and sustainable solutions into society









- Stakeholder Participants
- Discipline
- Goal, Shared Knowledge Academic Knowledge
- Thematic Umbrella Conventional Knowledge

Adapted from Wright Morton, L., S. D. Eigenbrode, and T. A. Martin. 2015. Architectures of adaptive integration in large collaborative projects. Ecology and Society 20(4):5.



- Within one academic discipline
- Disciplinary gal setting
- Development of new disciplinary knowledge

Multidisciplinary

- Multiple disciplines
- Multiple disciplinary goal setting under one thematic umbrella

Interdisciplinary

- Crosses disciplinary boundaries
- Development of integrated knowledge

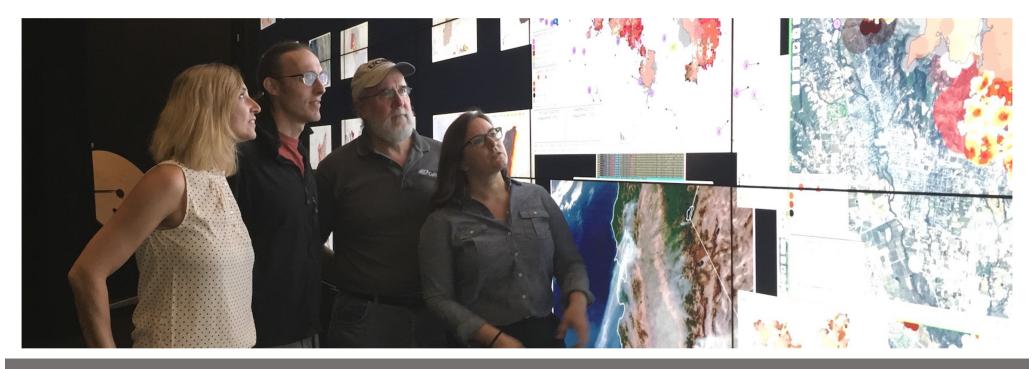


- Crosses disciplinary and sectorial boundaries
- Common goal setting
- Develops integrated knowledge for science and society
- Creates new paradigms



Translating Research into Impact

through Democratizing Access to Cyberinfrastructure





Three Main Components

Composable Workflows — Collaborative Innovation — Impact Network



- Develop methodologies and tools to enable collaborative workflow-driven data science
- Create solution architectures on top of big data and advanced computing platforms
- Push the boundaries of the computing continuum through composable systems and services



- Create collaborative pathways between UC San Diego and Los Alamos National Lab
- Accelerate the advancement of science and technology as a basis for responsible and proactive approaches to environmental challenges
- Leverage cutting-edge capabilities in scalable computing and diverse scientific expertise to foster solutionfocused, community-facing innovation



- Catalyze an impact network of students, researchers, practitioners, industry leaders, and public policy professionals committed to collaboratively engaging in research that is driven by specific and compelling societal problems and requires deep integration across disciplines and sectors to create solutions
- Provide participants with a foundational experience to position them for impact throughout their careers on the most challenging issues of our time.





CORE Institute Innovation Approach

Creating Breakthrough Technological Innovations for Complex Societal Challenges

Use-Inspired Problems

- 1) Context evaluation: Describe the system(s) within which the problem you are addressing exists and identify important decision-makers and vulnerable communities
- 2) Needs assessment: Clarify the needs of the people you want to help and ensure you are solving the right problem
- **3) Innovation pathways:** Sketch out ideas for data and science that could contribute to solving the problem and outline the expertise needed

Use-inspired & iterative

co-production of innovation

with users

CORE4 Building Blocks

Data & AI

Cutting-Edge Science & Engineering

Advanced Digital Infrastructure

Cutting-Edge Science & Integrated Workflows

Use-inspired & iterative

co-creation of solutions

with partners

Scalable Solutions

- 1) Sustainable partnership model: Implement solutions through a model that will allow for sustained use at scale
- 2) Continued iteration: Monitor performance and impact through user feedback and key metrics and be ready to adapt
- **3) Continued innovation:** Create mechanisms to ensure innovation is an ongoing process

From USEFUL



to USABLE



to USED at scale



İlkay Altıntaş, PhD (ialtintas@ucsd.edu)

UC San Diego ...
HALICIOĞLU DATA SCIENCE INSTITUTE

Translating Fire Research into Impact



Mission: Develop technologies with the fire management community driven by cutting-edge science and data

Vision: Enable tools that can have an impact at the scale of the environmental challenges we face today



wifire.ucsd.edu





Where are we headed at WIFIRE Lab?

- Wildfire Response: WIFIRE's Firemap platform in collaboration with CALOES and CAL FIRE through California's Fire Integrated Real-Time Intelligence System (FIRIS) and with partners in Colorado
- Beneficial Fire: WIFIRE's BurnPro3D platform for prescribed burn planning and implementation in collaboration with 3D fuel and fire modeling efforts at USGS, DOD, USFS, and LANL
- Data and Model Sharing: WIFIRE's Wildfire Technology Commons to develop standards, tools
 and techniques to share data and data-driven models to enable scientific workflows and AI
 innovation in collaboration with partners including NIST, CAL FIRE, and SDGE
- Immersive Visualization: Al-readiness of scientific data for new modes of teaching, training, decision-making, and public communication, including 3D outputs from vegetation modeling and fire science simulations and real-world information collected with cameras and sensors





Operational Products

FIREMAP

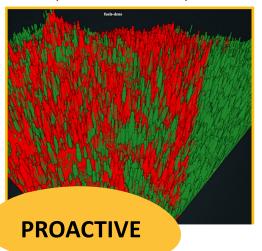
Firemap is currently being used by firefighters in Colorado, in collaboration with Intterra, and firefighters in California through the FIRIS program under the California Governor's Office of Emergency



Services and CALFIRE. FIRIS uses Firemap to provide realtime information on weather conditions and fire ignitions and to monitor and predict direction and speed of fire spread, as well as communities at risk. It has revolutionized initial attack response for the most dangerous fires across California.



In alignment with the nation's goal to increase fuel treatments to reduce wildfire risk, BurnPro3D is designed to support the preparation of burn plans as well as the implementation of prescribed



burns. The interface allows burn bosses to create and visualize high-resolution 3D fire simulations and compare fuel consumption and risk under different weather and ignition scenarios. It uses 3D FastFuels data developed by the US Forest Service and the QUIC-Fire coupled fire/atmosphere model developed at Los Alamos National Lab.











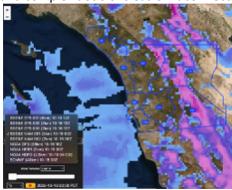






Wildfire Science and Technology Commons

The Commons enables the development of foundational AI techniques to fuse and learn from data and to make scientific models interpretable and complex decisions easier. It connects next-generation data and





models for anyone interested in developing solutions. For example, it enables an integrated fire weather intelligence platform focused on reducing risk related to power lines for Southern California. A new phase of development was recently supported through congressionally directed spending proposed by California Sen. Padilla, Rep. Vargas, and Rep. Jacobs.

Wildfire and Landscape Resilience Data Hub

The Data Hub is a federated data ecosystem for California's Wildfire and Forest Resilience Task Force, providing a "single view" over existing data to fulfill the reporting requirements for California's



Million Acre Strategy to treat 1 million forested acres per year to reduce wildfire risk. It will provide public, open, and fair access to data, analytic tools, and customizable reports via the Data Hub explorer web viewer, as well as access to data through APIs.





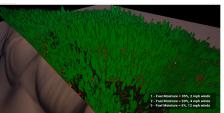


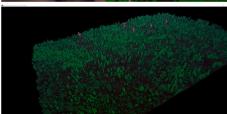


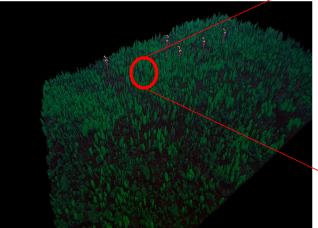


Immersive Forest for **Multimodal Communication**











Terrestrial LiDAR contextualized within Aerial scan







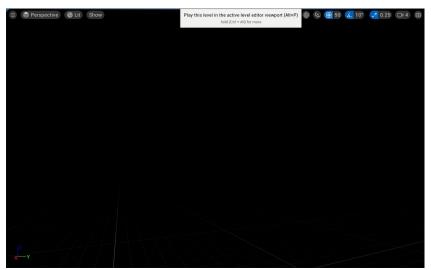
AI in Science Communication

Visualization of multiple terrestrial LiDAR scans in the Immersive Forest prototype



Immersive Al-integrated visualization of scientific data and simulations for training, decision making, and public communication.

Animations by: Isaac Nealey (left, bottom), Ivannia Gomez (top)









Additional Grants Fueling R&D



Evaluation of satellite-based fire detection and fire radiative power applications



Ground sensing and in-situ edge computing for monitoring and decision-making



Open fire models to predict wildfire spread over 3-5 days



Workflows for DOD prescribed fire managers participating in the National Innovation Landscapes Network



Prescribed fire planning and monitoring tools and workforce training for California agencies



Multi-modal data to improve characterization of fuels at large spatial extents and fine spatial scales



Immersive
visualization of
scientific data for
new modes of
training,
decision-making
and
communication



This type of work needs the CORE4 building blocks.

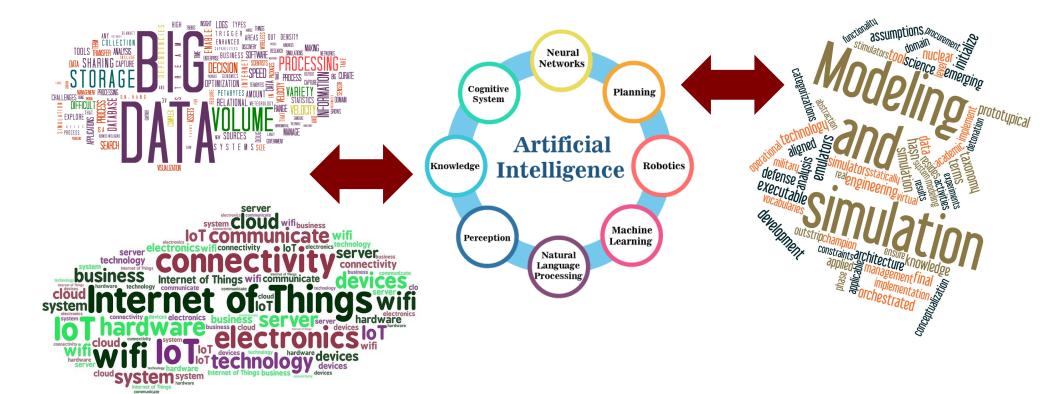


CORE4 Building Blocks

Data & Al	Cutting-Edge Science & Engineering
Advanced Digital Infrastructure	Integrated Workflows



Al-Integrated Applications at the Digital Continuum













Al in Science and Research 2023





Al in Science Readiness "not just science + Al methods"

Data federation and hubs

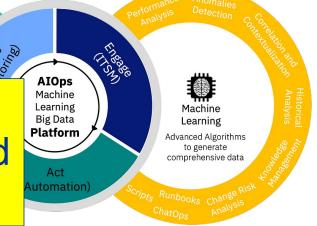
Data quality and volume

Knowledge manager

Benchmarks

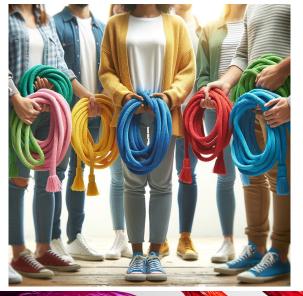
Workflow managem

Requires a full team Scalability up and do and enabling integrated data platforms



- Software integration and engineering
- Dev ops (also called AI ops and data ops)
- Interpretability and explainability
- Workforce training and culture/incentive building



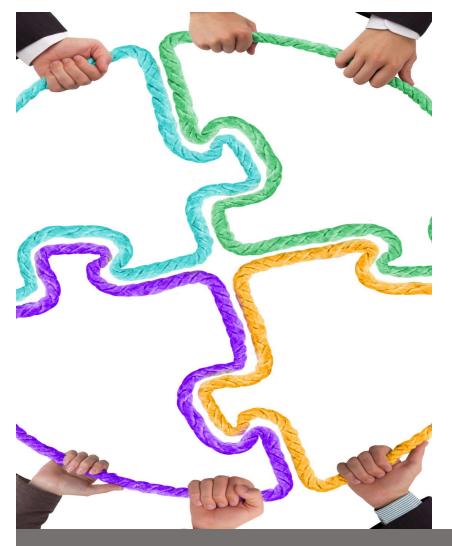


Systems should enable seamless integration of Al-integrated application workflows by teams!





UCSanDiego HALICIOĞLU DATA SCIENCE INSTITUTE

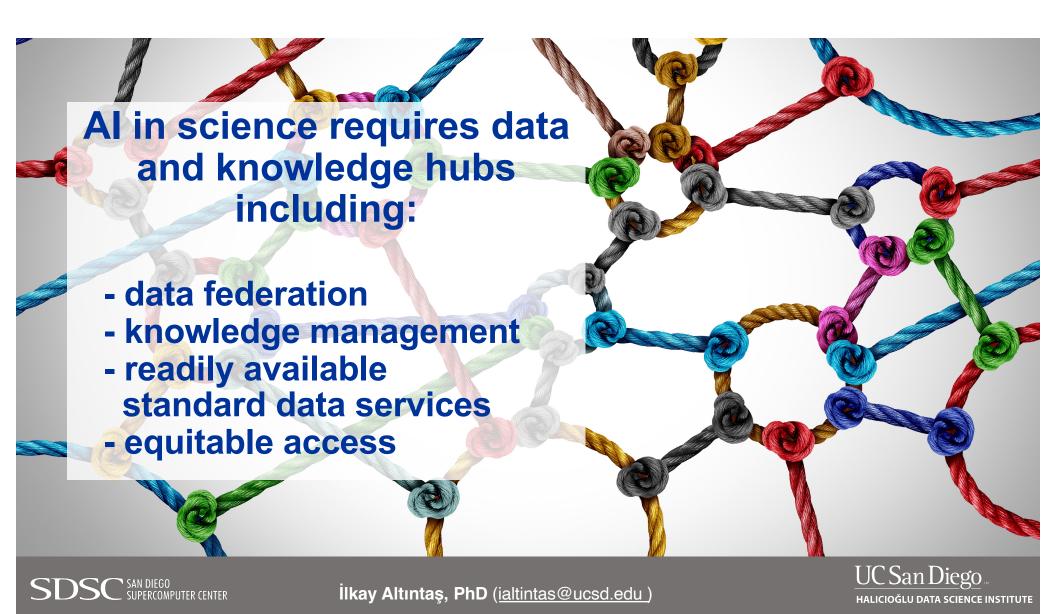


Workflow integration requires a digital continuum composed through:

- system federation
- reusable capability services
- solutions integrating services







Integration requirements...



Dynamic composability matters.

Systems and services are useful if groups can integrate them into applications.





Tools that enhance teamwork and use need to be coupled with responsible AI systems.





Dynamic composability matters.

COMPOSABLE SERVICES

e.g., model and data archives, learning and analytics, simulation, training

RESOURCE MANAGEMENT

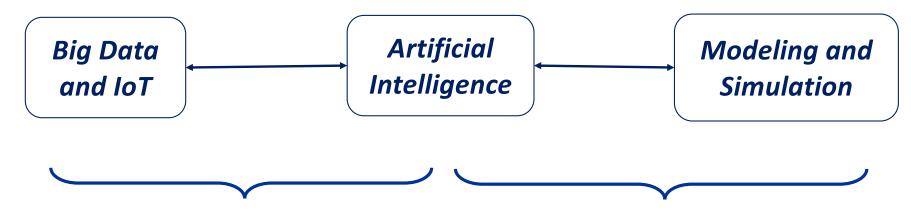
e.g., container orchestration, optimization

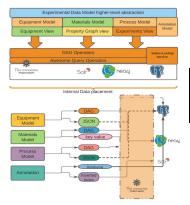
COMPOSABLE SYSTEMS

e.g., GPU, CPU, Big Data, quantum, neuromorphic, SDN, storage









Big Data





Capacity



xPU → GPU, CPU, TPU, IPU, QPU,...









Cloud, HPC, Storage



Some Composable Systems









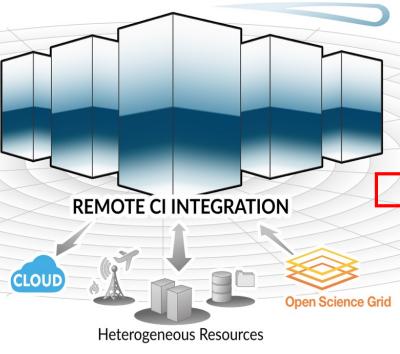
EXPANSE COMPUTING WITHOUT BOUNDARIES 5 PETAFLOP/S HPC and DATA RESOURCE

HPC RESOURCE

13 Scalable Compute Units 728 Standard Compute Nodes 52 GPU Nodes: 208 GPUs 4 Large Memory Nodes

DATA CENTRIC ARCHITECTURE

12PB Perf. Storage: 140GB/s, 200k IOPS
Fast I/O Node-Local NVMe Storage
7PB Ceph Object Storage
High-Performance R&E Networking



LONG-TAIL SCIENCE

Multi-Messenger Astronomy Genomics

Earth Science Social Science

INNOVATIVE OPERATIONS

Composable Systems

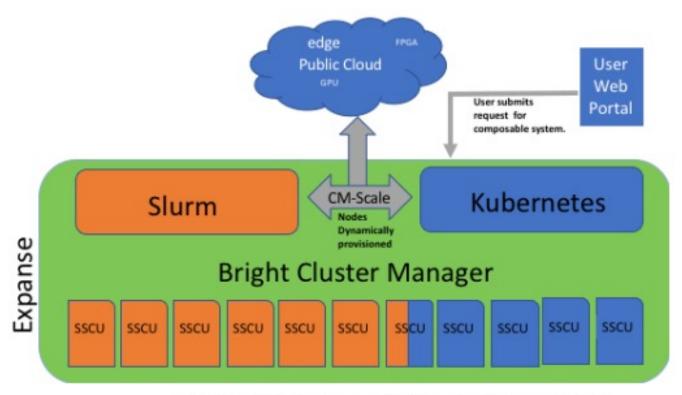
High-Throughput Computing Science Gateways

Interactive Computing
Containerized Computing

Cloud Bursting



UC San Diego ...
HALICIOĞLU DATA SCIENCE INSTITUTE



Expanse Composable Systems Framework





National Research Platform



https://nationalresearchplatform.org/







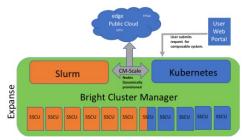
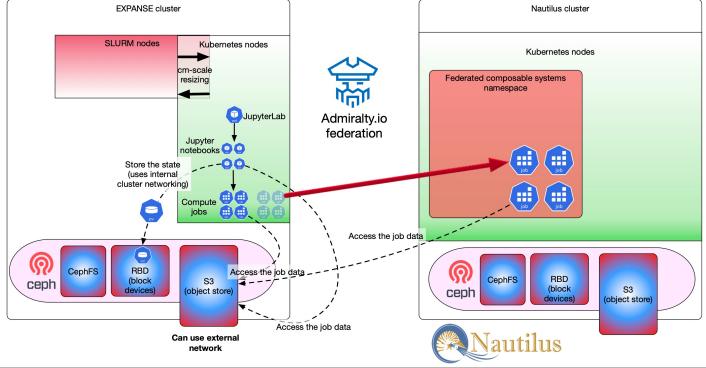


Figure 5.1 Expanse Composable Systems Framework

First composable cluster is federated!

EXPANSE (Enthalpy) + CHASE-CI (Nautilus)









HPC/Cloud

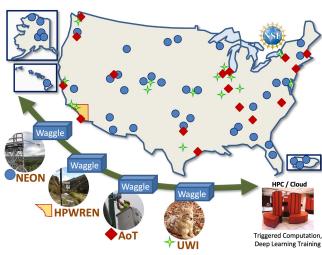
Al@Edge and the Digital Continuum

Slide Source: Pete Beckman, ANL



























neon









LINCOLN PARK ZOO.



Education & Training









long-term archive MANAGEMEN data reuse services active data repositories, networks, knowledge

Systems and services are only useful if the groups can integrate them into applications.

WORKFLOW MANAGEMENT

e.g., application integration, coordination, optimization, communication, reporting

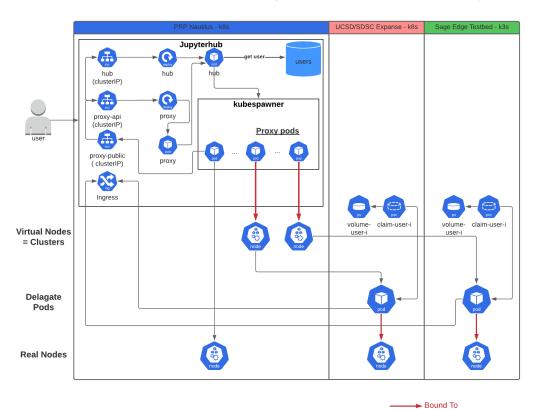
COMPOSABLE SERVICES

RESOURCE MANAGEMENT

COMPOSABLE SYSTEMS



Integration of NSF EXPANSE, NRP and Sage A Composable System Deployment of JupyterHub





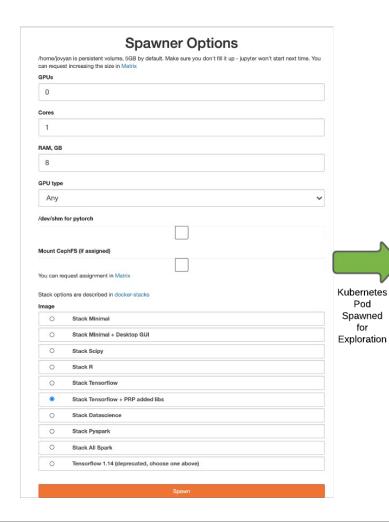
PRP Nautilus

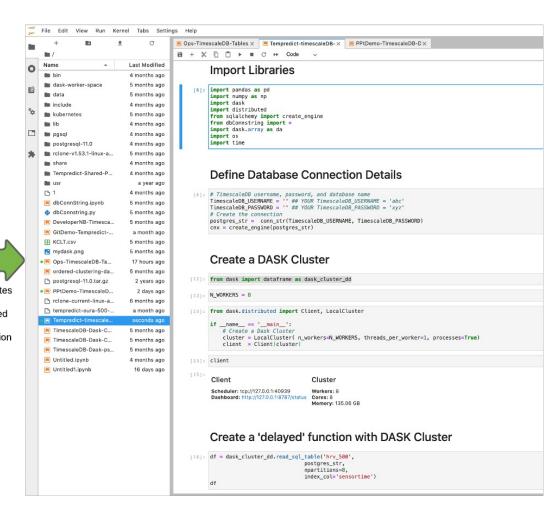
- A user can easily be provided the right environment for developing their AI Edge Application
- I. Altintas et al., "Towards a Dynamic Composability Approach for using Heterogeneous Systems in Remote Sensing," 2022 IEEE e-Science doi: 10.1109/eScience55777.2022.00047

deploy on the Edge





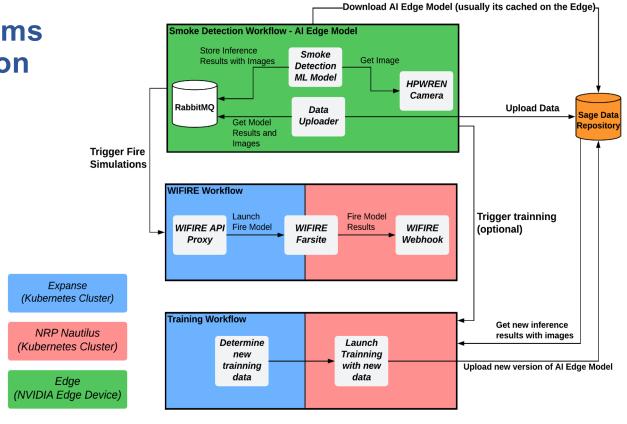






Fire Simulations using Composable Systems and Edge Smoke Detection

- Three workflows
 - Smoke Sage Edge App
 - Fire simulator
 - Al Training
- Both the fire simulator and training workflows are can be run on Expanse or Nautilus through the federation layer



I. Altintas et al., "Towards a Dynamic Composability Approach for using Heterogeneous Systems in Remote Sensing," 2022 IEEE e-Science doi: 10.1109/eScience55777.2022.00047





The king's horses ran simulations fast,
While the king's men studied data amassed.
Yet despite their efforts, and all they'd apply,
True AI readiness still passed them by.

İlkay Altıntaş, PhD (ialtintas@ucsd.edu)



accuracy, privacy, explainability, ethics

e.g.,

TEAM SCIENCE

USE-INSPIRED INTERFACES

e.g., for science, education and scalable practice

Tools that enhance teamwork and use need to be coupled with responsible AI systems.





ethics, equity accuracy, privacy, explainability, RESPONSIBIL

e.g.,

REPRODUCIBILITY

TEAM SCIENCE

long-term archives LIFECYCLE MANAGEMEN data reuse services active data repositories, knowledge networks, DATA

USE-INSPIRED INTERFACES

e.g., for science, education and scalable practice

WORKFLOW MANAGEMENT

e.g., application integration, coordination, optimization, communication, reporting

COMPOSABLE SERVICES

e.g., model and data archives, learning and analytics, simulation, training

RESOURCE MANAGEMENT

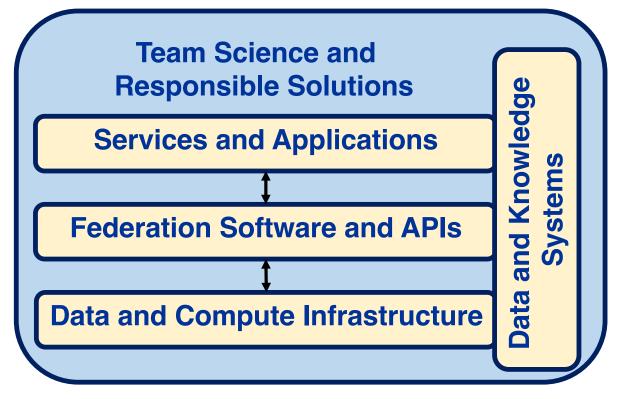
e.g., container orchestration, optimization

COMPOSABLE SYSTEMS

e.g., GPU, CPU, Big Data, quantum, neuromorphic, SDN, storage



Use-Inspired Composability from Systems to Services



- User-centered design and experience
- Improved FAIR data capacity
- Capability-based integration
- Create plug and play microservices
- Run across many systems
- Dynamically measure, manage and provision resources



Schmidt AI in Science Postdoc Research Schmidt Sciences



Computational microscopy of respiratory viruses in aerosols

> Exploring different models to simulate and visualize the behavior of viruses in the respiratory tract

The relationship between life span of the plant roots microscopy data and wildfire

> Deep learning model to estimate life span

AI-Powered analysis of molecular simulations

> High-affinity generative model for target proteins

Small coronary artery calcium detectability

> Deep learning model to segment and visualize chambers of the heart

Data-driven development of neural-network potentials from quantum chemistry data

> ML model to be used as a surrogate for expensive highlevel chemistry calculations

Drug resistance evolution in HIV patients

> Leverages machine learning system for heterogeneous cryo-EM reconstruction of proteins and protein complexes from single-particle cryo-EM data

Earth system modelling

Deep learning model to use data extracted from ECMWF to calibrate earth systems simulation

Brain activity of diving seals reveals short sleep cycles at depth

> Linear regression models to assess the impact of age, recording location and design

Bathymetry from space

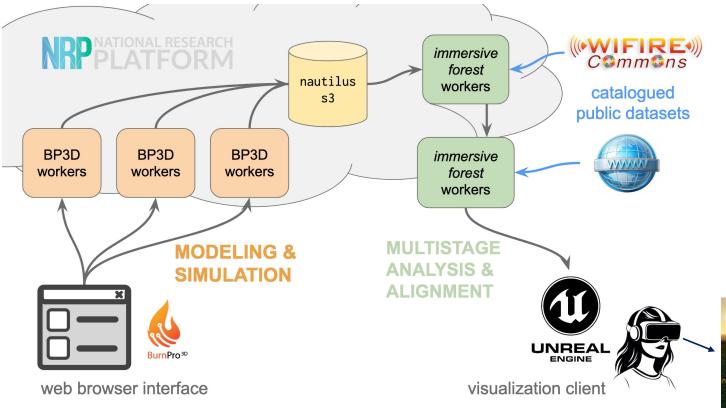
Machine learning model to understand small-scale ocean dynamics

The effect mutations implicated in autism can have in protein oscillation

Deep learning model to predict the oscillation of protein in cellcell communication







Current prototype capabilities

- Terrestrial LiDAR contextualized with aerial LiDAR for VR
- Georeferenced panoramic projection of terrestrial LiDAR for mobile
- Watch and interact with fire simulations in 3D under a variety of weather conditions
- Move through multiple LiDAR scans across the landscape to compare pre- and post-burn vegetation in 3D

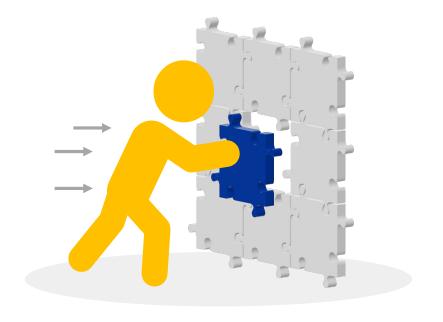


3D Immersive Forest using NRP





What do we do about the data gaps?



http://www.nationaldataplatform.org























http://www.nationaldataplatform.org



Award abstract: https://www.nsf.gov/awardsearch/showAward?AWD ID=2333609



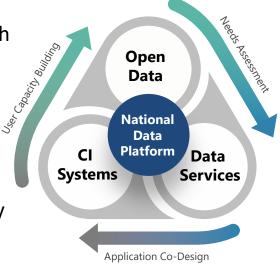


National Data Platform Pilot (NDP): Services for Equitable Open Access to Data

A **federated** and **extensible** data ecosystem to promote collaboration, innovation and equitable use of data using existing and future national cyberinfrastructure (CI) capabilities.

 A broad data ecosystem to enable data-enabled and AI-integrated research and education workflows

- Facilitates data registration, discovery and usage through a centralized hub
- Enhances distributed CI capabilities through distributed points of presence
- Cultivates resources for classroom education and data challenges
- Assists research and learning through personalized workspaces
- Partnership pathways to foster scientific discovery, decision-making, policy formation and societal impact
 - Focus areas: Wildfire, climate, earthquake and food security, among others





Addressing Open Questions for Equitable Open Access

Foundational
Abstractions and
Services

- What are the foundational data abstractions and services that can serve as multipurpose and expandable building blocks for data-driven and Al-integrated application patterns?
- How can everyone effectively access and utilize these abstractions and services?

Equitable and Open CI Use

- How can such foundational data abstractions and services be developed and deployed on top of existing production-ready CI, including storage and the edge-to-HPC continuum?
- How can we ensure equity of data access and use across distributed CI?

Needs, Requirements and Challenges

- What are the requirements and challenges for governance of open science, open data and open CI?
- What are the required guardrails for protecting privacy, civil rights and civil liberties that will ensure a more equitable use of data systems and services?



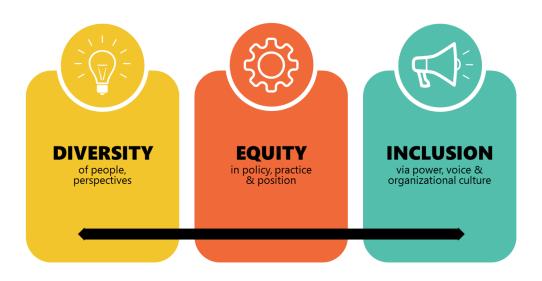
Diversity is a fact.

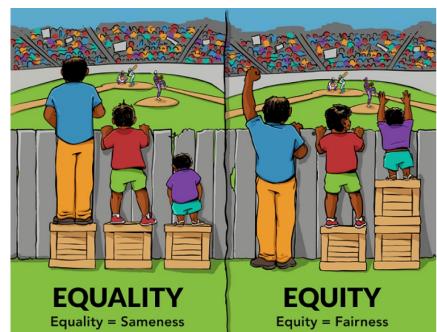
Equity is a choice.

Inclusion is an action.

Belonging is an outcome.

- Arthur Chan





Equality promotes fairness and

justice by giving everyone the same thing.

BUT, it can only work if everyone starts from the same place. In this example, equality only works if everyone is the same height. Equity is about making sure people get access to the same opportunities.

Sometimes our differences or history can create barriers to participation, so we must FIRST ensure EQUITY before we can enjoy equality.

Source: Angus Maguire for the Interaction Institute for Social Change http://interactioninstitute.org/illustrating





EXECUTIVE OFFICE OF THE PRESIDENT OFFICE OF SCIENCE AND TECHNOLOGY POLICY WASHINGTON, D.C. 20502

August 25, 2022

MEMORANDUM FOR THE HEADS OF EXECUTIVE DEPARTMENTS AND AGENCIES

Dr. Alondra Nelson Andre Nulson

Deputy Assistant to the President and Deputy Director for Science and Society Performing the Duties of Director

Office of Science and Technology Policy (OSTP)

SUBJECT: Ensuring Free, Immediate, and Equitable Access to Federally Funded Research

This memorandum provides policy guidance to federal agencies with research and development expenditures on updating their public access policies. In accordance with this memorandum OSTP recommends that federal agencies, to the extent consistent with applicable law:

- 1. Update their public access policies as soon as possible, and no later than December 31st, 2025, to make publications and their supporting data resulting from federally funded research publicly accessible without an embargo on their free and public release;
- 2. Establish transparent procedures that ensure scientific and research integrity is maintained in public access policies; and,
- 3. Coordinate with OSTP to ensure equitable delivery of federally funded research results





Empowering citizens & strengthening accountability

- Increases citizen engagement

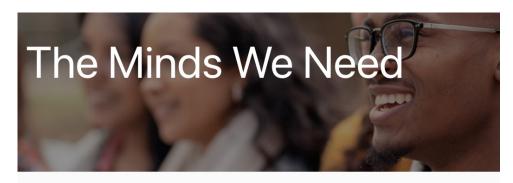


Innovation & efficiency in government agencies

- Inter-agency collaboration

for the economy

services for the entire economy



Inclusion, Innovation, and Competitiveness | Strengthening Our National Broadband Initiative | Investing in Research and Education Infrastructure | Contributors | Toolkit | Endorsements

Inclusion, Innovation, and Competitiveness

We are at a crossroads.

https://mindsweneed.org

Toward Democratizing Access to Facilities Data: A Framework for Intelligent Data Discovery and Delivery

Data collected by large-scale instruments, observatories, and sensor networks Data contected by tall gescale installments, observations, and serian networns (i.e., science facilities) are key enablers of scientific discoveries in many disciplines However, ensuring that these data can be accessed, integrated, and analyzed in a democratized and timely manner remains a challenge. In this article, we explore how state-of-the-art techniques for data discovery and access can be adapted to facilitate data and develop a conceptual framework for intelligent data access an



Democratizing Computation and Data to Bridge Digital Divides and Increase Access to Science for Underrepresented Communities

October 3 2021 NSF OAC 2127450

Democratization of Cl and Data Access



UC San Diego HALICIOĞLU DATA SCIENCE INSTITUTE

Architecting for Collective Data-Integrated Impact

- Involve diverse users in architecting
- Identify access, use, expertise and education gaps
- Improve the experience of working with data
- Connect data to knowledge systems and services
- Create an ecosystem approach to capacity building
- Incubate use-inspired solutions to scale
- Explore new models of allocation
- Develop and teach models of sustainability and scale

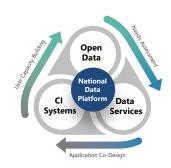






Our Use-Inspired Approach

Solving data gaps one workflow template at a time...



Identify Gaps

- Community advisory board
- External community integration plan
- Needs assessments
- Co-design workshops
- Expansion prototypes



Incubate, Innovate and Educate

Use-Inspired Workflows and Interfaces

Data and Knowledge Management



Composable Systems and Platforms

Sustainable and Scalable Use

- Distributed in nature
- Composition as a principle
- Hub-centric services as connection backbone
- Integrates in education systems

Collaboration, Incubation, Allocation and Partnership Models



UC San Diego HALICIOĞLU DATA SCIENCE INSTITUTE







Centralized portal for discovery, access and use workspaces for research and education



A scalable **platform** for using, developing and deploying services and application workflows at **distributed** points of presence



Current 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 NDP Federation** Central discovery & access Scalable platform for customizable & workspace for research & education composable service stacks NDP Portal NDP POP NDP POP Search, **Education Hub** Discovery, NDP POP NDP POP Access & Use **Coordinated Crosscutting Services** Authentication & Authorization, Accounting, Orchestration, Monitoring and Logging, Catalog, Workspace, Notebooks, Data Democratization, Workflow Integration, etc. Pelican/OSDF Cloudbank **Nautilus ACCESS** (\ldots) SAGE jupyter Data CI and Delivery Services Containers **HPC** Cloud Composable CI



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

NDP Hub

Search,
Discovery,
Access & Use

- 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
- Al 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

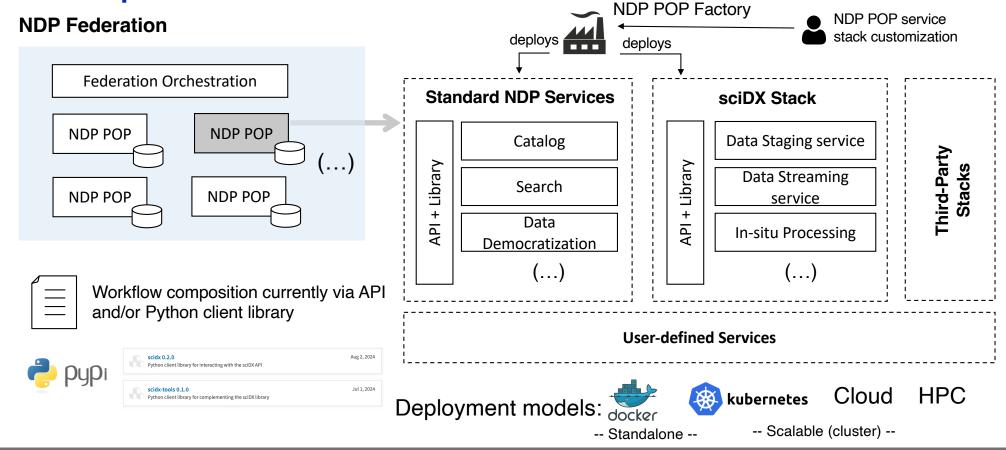
Planned Future Work

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





NDP POP: Distributed Points of Presence with Customizable, Composable Service Stacks





Typical NDP Workflow with Composable Capabilities







Data sources:

NAIRR datasets, repositories, instruments, sensors, facilities, etc. Data adaptors: Customizable metadata/data ingestion; heterogeneous data sources & types

Data Acquisition

Discovery

(Meta) data

Curation.

Registration,

Indexing,

- Metadata catalogExtensible search
- engine
- Recommendation services

Staging, Insitu Data Analytics

Data

Access.

Data

- High-performance data staging & in-situ processing
- Data access optimization (caching, pre-fetching, recommendation)
- Leverage Pelican data origins; data federation

- Data Processing, Data-driven, AI/ML-based Workflows,
 - **⇒**
- Generation, Curation, Sharing, Archival

Product

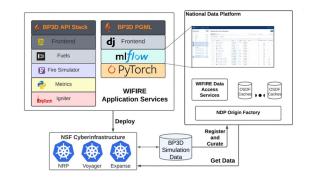
- Services Gateway (support notebooks)
- Leverage NAIRR/NSF
 ACCESS/Cloud resources.
- Data-driven stream processing
- Data product & re-streaming
- Archival support (including using Pelican)

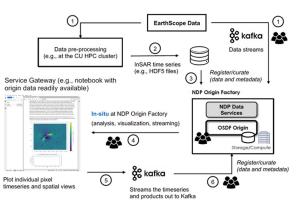


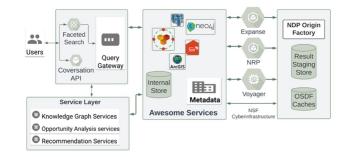


Case Studies for Generalizable Workflows

- Representative examples of important patterns that exist in science today for working with
 - O large datasets
 - O 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.















Planned Extensions for NAIRR (September 2024 – August 2025)



NAIRR Data Resource Catalog

- Ingestion Process for NAIRR Data
- FAIR NAIRR Catalog
- Conversational Search Interfaces

NAIRR CloudBank Research Workflows

- Provisioning and Accounting
- CloudBank Workflow Deployment
- Collaborate with NAIRR science pilots

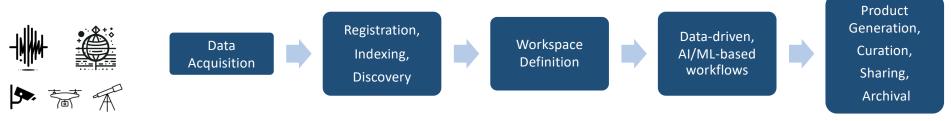
NAIRR Classroom Workflows

- NAIRR Educator Workflows
- NAIRR Student Workflows
- Community Engagement

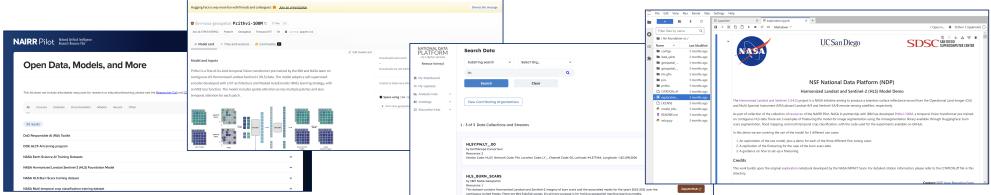




Example NDP-NAIRR AI in Science Workflow



- 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 Al/ML workflow is supported by Al 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.





UC San Diego ...
HALICIOĞLU DATA SCIENCE INSTITUTE

NDP Hub Functionality

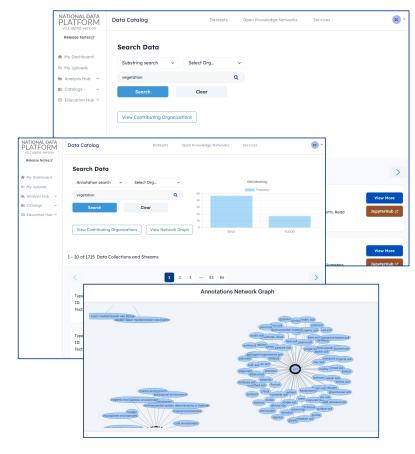
September 2024 Release

60





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.

Future Work Post-September 2024 to include NAIRR Data Resources:

- 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.

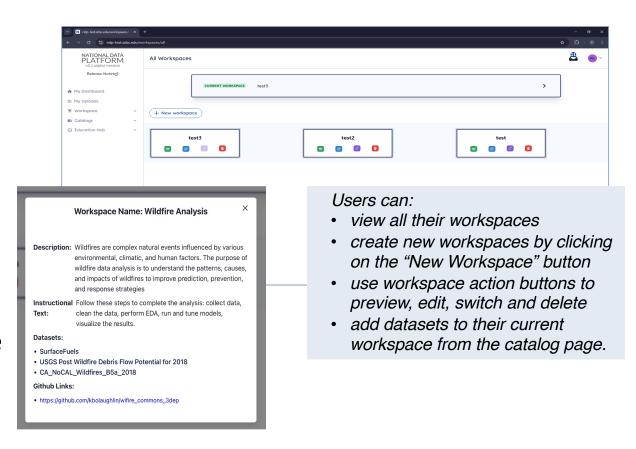




NDP Workspaces (Version 1 – September 2024)

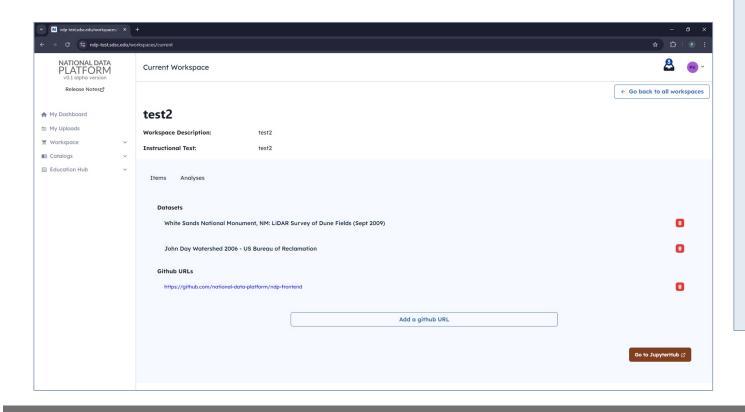
Goal: Craft persistent and customizable workspaces with datasets and services to launch into a sandbox

- 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





Current Workspaces



Users can:

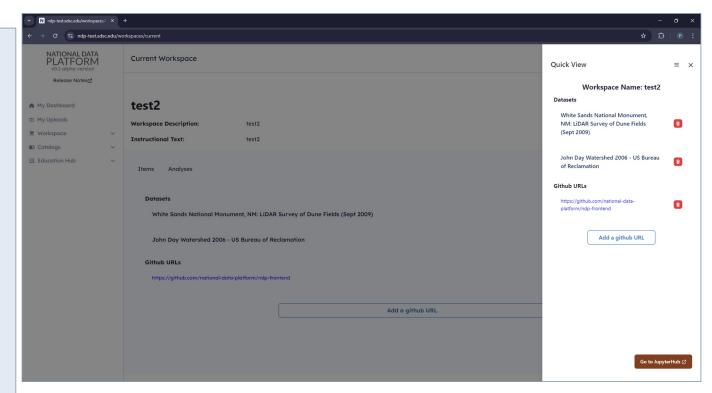
- view their current workspace along with all the details and resources
- delete datasets and add/delete github links from their current workspace
- navigate to
 JupyterHub using
 the button at the
 bottom



Current Workspace Quick View Drawer

Users can:

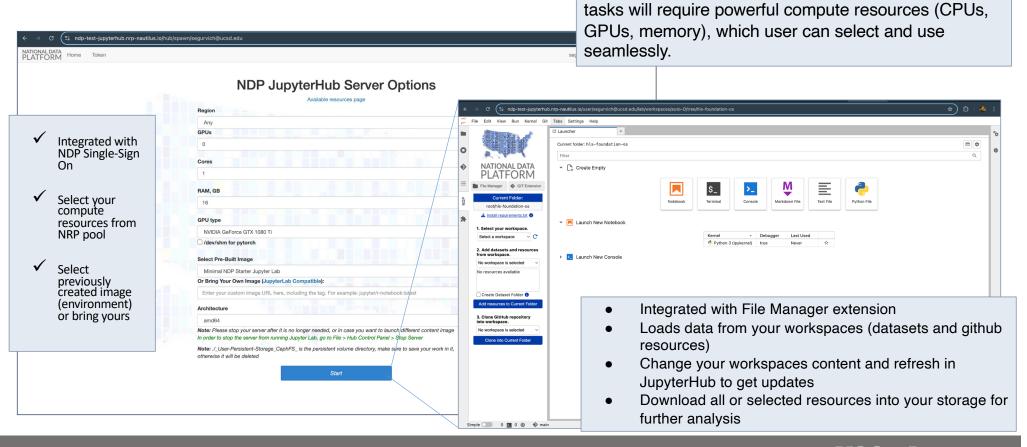
- see a quick view of their current workspace along with all the details and resources
- enter this view by clicking the drawer button on the top right on the navbar next to their avatar
- delete datasets and add/delete github links from their current workspace
- navigate to JupyterHub using the button at the bottom







NDP JupyterHub (Sandbox)

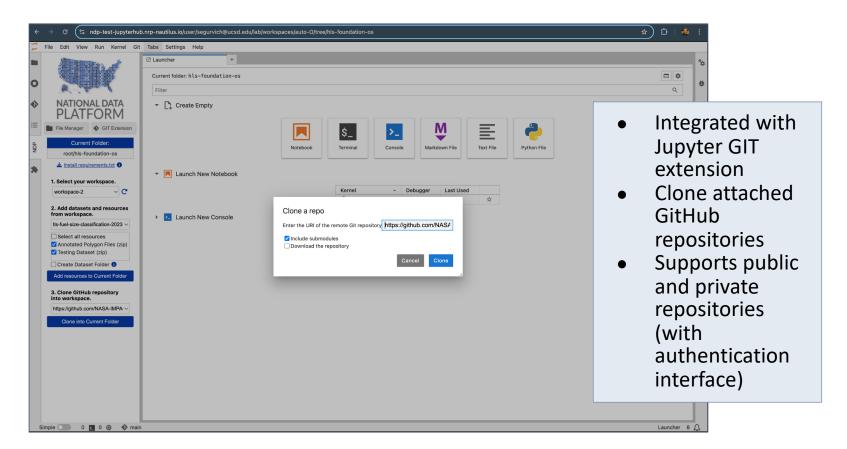




A compute environment for data analysis, machine

learning training or any other computational tasks, built on top of NRP (Nautilus) cluster. Different datasets and

NDP JupyterLab Extension





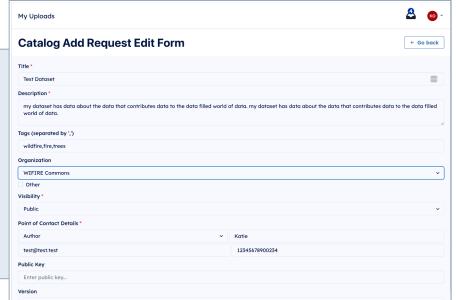
NDP Catalog Addition

Goal: Users can add dataset references to either NDP centralized catalog or POP-specific catalog



Curated Public Catalog Add Request:

- Provide all metadata and data access information
- Designated data approvers evaluate dataset quality
- Add or reject datasets for access to community





NDP Data Challenges for students and researchers

NDP Education Gateway to provide participants access

The challenge questions will require using data and models in an environment that requires computing and huge data stores, which would typically be

unavailable to a student or

researcher without the NDP

Education Gateway.

Designed to ensure that we are

services for equitable education

developing broadly accessible

and community building.

Three Co-Design Workshops

to the NDP data ecosystem

Each will include a breakout session to develop a data challenge question specific to large data (W1); streaming data (W2); and graph data (W3).

Education and capacity building through data challenges



Data challenge toolkits will be developed after each data challenge so that other institutions can easily design their own data challenges to be run through the NDP Education Gateway.

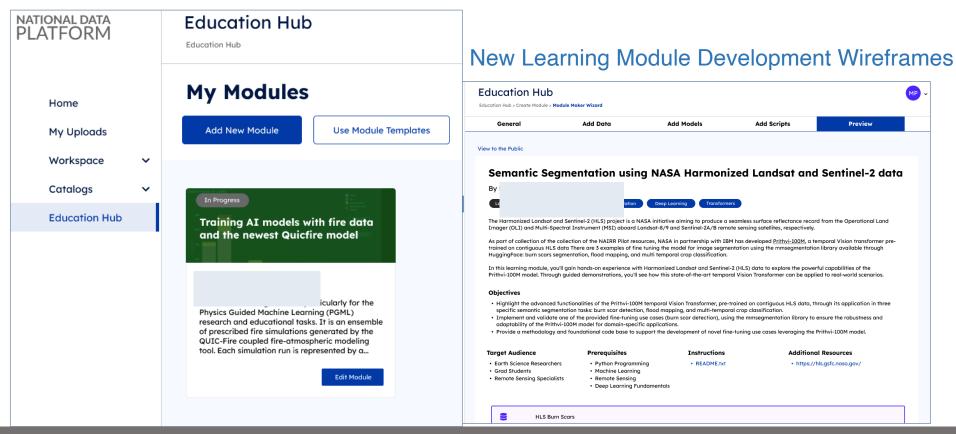








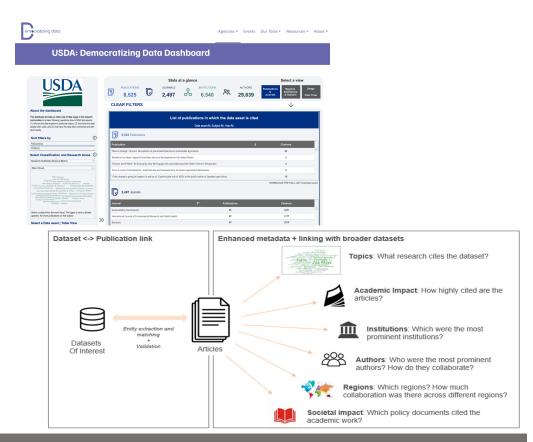
NDP Education Hub (Version 1 – September 2024)





Democratizing Data (DD) Service

- Composable DD service for search and discovery within NDP requires a detailed and structured approach
 - Extracting and utilizing publication metadata using multiple corpora (e.g., Scopus, OpenAlex, PubMed)
 - Integrating NAIRR datasets starting with USDA, NIH, NASA, and NOAA.
 - Exploring a generalized approach to support the integration to other corpora and Al-ready NAIRR datasets
- Leveraging https://democratizingdata.ai/





NDP PoP Examples & Documentation

September 2024 Release

71





Science Data Exchanges (sciDX) Services: Data Staging and Streaming Services

Science Data Exchange (sciDX): Customizable software stack for in-situ data access & processing

Data Staging Service

- In-situ (close to the data) data processing and access
- High-performance in-memory processing
- Server-side data transformations (e.g., subsetting, reduction, user-defined analysis, etc.)
- Caching/sharing of data, query results, and data products with user and group isolation

Data Streaming Service

- Streams registration, curation/archival for discovery and access
- User-defined operations on streaming data (semantically specialized abstractions)
- Combine streaming data with archived/playback data
- Mechanism for online data product generation (i.e., new data streams

In-situ AI workflow execution runtime (on staged and streaming data)

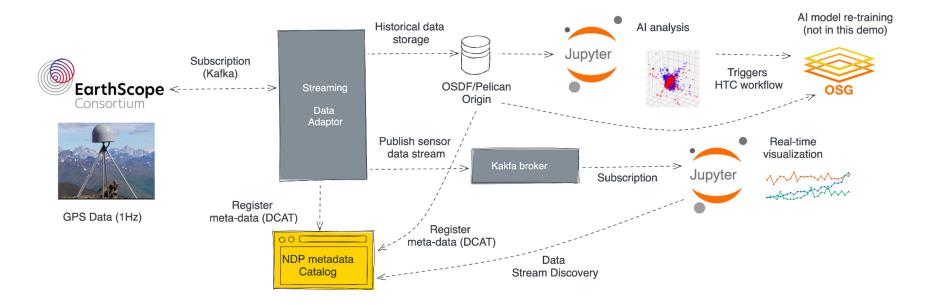




Example 1: EarthScope data streaming/analysis enabled by NDP POP

Real-time high-precision GNSS stations

- Al analysis: anomaly detection from archived data from OSDF/Pelican
- Real-time data visualization

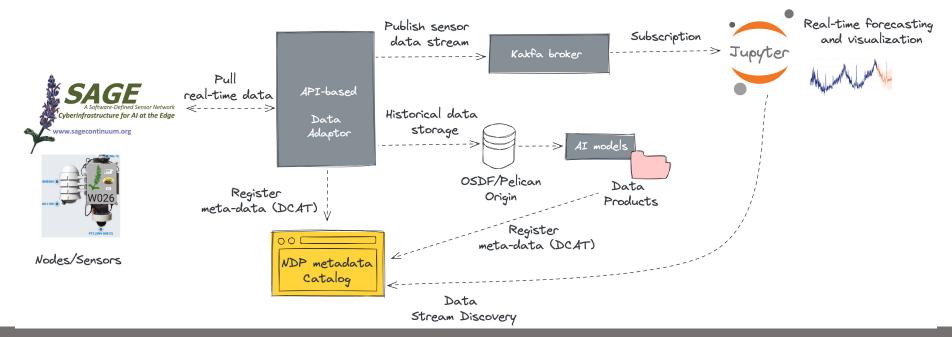




Example 2: SAGE data streaming/analysis enabled by NDP

SAGE data streams

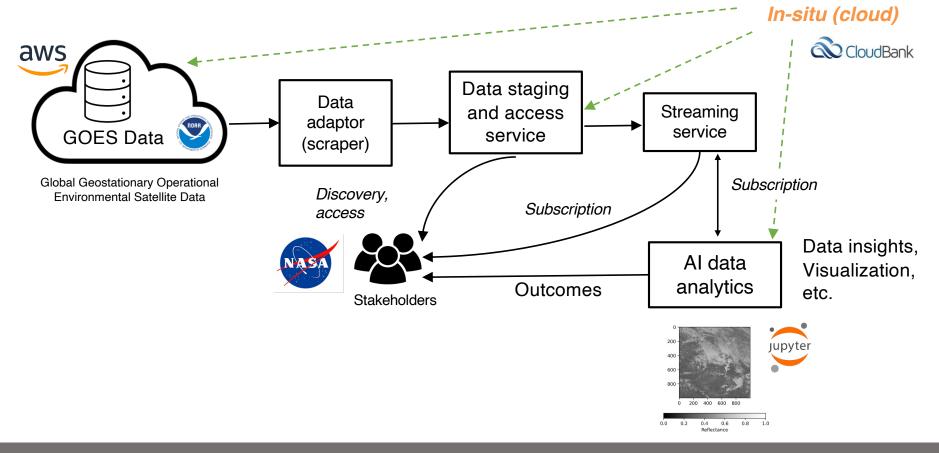
- Real-time data visualization (temperature)
- Time series forecasting (proof-of-concept)





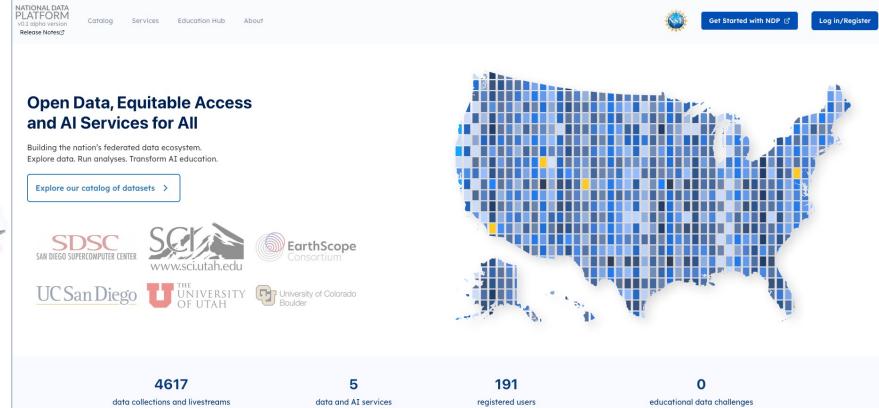


Example 3: Fire Detection using NASA GOES Satellite Data (under development)





Any questions? Contact ndp@sdsc.edu







To sum up...

Emerging new applications require integrated AI in dynamically composed workflows.

Artwork: **Jen Stark, Cosmographic, 2014**, acid-free paper, holographic paper, glue, wood, acrylic paint, 34 x 37 x 4 in.



Complexity comes at a cost

- Composable systems is not a turnkey functionality
- Requires collaboration with and between infrastructure providers

Convergence research helps

- End-to-end data pipelines need to be defined for each application along with microservice execution
- Use-inspired design and translational CS helps to focus the effort

