# Data Infrastructure in the AI Era

Keynote Talk for the International Symposium on Grids and Clouds
March 19, 2025 – Taipei, Taiwan

**İlkay ALTINTAŞ, Ph.D.**

**University of California, San Diego**

Chief Data Science Officer & Division Director, Cyberinfrastructure and Convergence Research, **San Diego Supercomputer Center**
Founding Fellow, **Halıcıoğlu Data Science Institute**
Founding Director**, Workflows for Data Science Center of Excellence**
Founding Director, **WIFIRE Lab**

Joint Faculty Appointee, **Los Alamos National Laboratory**

SDSC SAN DIEGO SUPERCOMPUTER CENTER

UC San Diego ™
HALICIOĞLU DATA SCIENCE INSTITUTE

SCHOOL OF COMPUTING, INFORMATION AND DATA SCIENCES

UC San Diego

School of Computing, Information and Data Sciences

https://scids.ucsd.edu/

SAN DIEGO SUPERCOMPUTER CENTER
at the UNIVERSITY OF CALIFORNIA
COMPUTING WITHOUT

ABOUT SDSC    SERVICES    SUPPORT    RESEARCH & DEVELOPMENT    ED

Materials Science Researchers Double Up on SDSC, PSC Supercomputers to Discover New Details about TMDs

Supercomputer simulations provide a better understanding of two-dimensional layered materials showing promise for a variety of applications – from flexible electronics and spintronics to opto and memory devices.

READ MORE

Innovate,

FOR
UC/UCSD Researchers

FOR
National HPC Users

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

UPCOMING EVENTS

OCT 2    2:00 pm - 3:00 pm
Some new results for streami

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

View Calendar

About    People    Research    Graduate    Undergraduate    Industry    Connect With Us    Get Invo

NEWS

Pioneering Data Science for a Data-Driven Future

JULY 18, 2023 — KALEIGH O'MERRY

Tweets from @HDSIUCSD

Nothing to see here - yet

When they Tweet, their Tweets will show up here.

View on Twitter

w Does
atGPT

NEWS

How Does ChatGPT Work? – Event

cience.ucsd.edu/

SDSC SAN DIEGO SUPERCOMPUTER CENTER

**İlkay Altıntaş, PhD** (ialtintas@ucsd.edu )

UC San Diego
HALICIOĞLU DATA SCIENCE INSTITUTE

# Cyberinfrastructure and Convergence Research Division @SDSC

### Translating cyberinfrastructure research for impact at scale

## CI Methods and Systems

- "Big" Data and Knowledge Systems
- Computational Data Science
- Machine Learning and AI
- Advanced Computing

## Convergence Research

- Collaborative Problem Solving
- Use-inspired Design
- Sustainable and Scalable Solutions

## Experiential and Classroom Education

WCRP
# GRAND CHALLENGES

MELTING ICE AND GLOBAL CONSEQUENCES

CLOUDS, CIRCULATION AND CLIMATE SENSITIVITY

REGIONAL SEA LEVEL CHANGE AND COASTAL IMPACTS

WATER FOR THE FOOD BASKETS OF THE WORLD

WEATHER AND CLIMATE EXTREMES

CARBON FEEDBACKS IN THE CLIMATE SYSTEM

NEAR-TERM CLIMATE PREDICTION

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

## CO₂ emissions per capita vs GDP per capita

Our World in Data

Per capita consumption-based CO₂ emissions

CO₂ emissions are too high

Energy poverty

To end climate change the long-run goal is that net-emissions decline to zero.

Data for 2017: Global Carbon Project, UN Population, and World Bank.
OurWorldinData.org – Research and data to make progress against the world's largest problems.
Licensed under CC-BY by the author Max Roser.

https://ourworldindata.org/worlds-energy-problem

$10 trillion+
spent in global response to COVID-19

216
countries, areas or territories with cases

165+
COVID-19 vaccines being developed globally

https://www.10xgenomics.com/research-areas/infectious-disease

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

SDSC SAN DIEGO SUPERCOMPUTER CENTER

**İlkay Altıntaş, PhD** (ialtintas@ucsd.edu )

UC San Diego
HALICIOĞLU DATA SCIENCE INSTITUTE

**Convergence research is:**

**driven by a specific and compelling societal problem**

**and**

**works towards integrating innovative and sustainable solutions into society**



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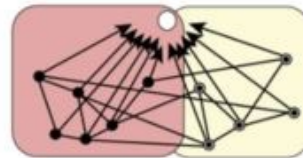- Development of new disciplinary knowledge

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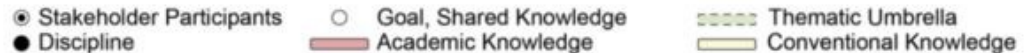- Multiple disciplinary goal setting under one thematic umbrella

**Interdisciplinary**
- Crosses disciplinary boundaries
- Development of integrated knowledge

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

◉ Stakeholder Participants   ○ Goal, Shared Knowledge   ▭ Thematic Umbrella
● Discipline                 ▭ Academic Knowledge        ▭ 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.

İlkay Altıntaş, PhD (ialtintas@ucsd.edu )

# Translating Research into Impact
## through Democratizing Access to Cyberinfrastructure



İlkay Altıntaş, PhD (ialtintas@ucsd.edu )

# Three Main Components

**Composable Workflows** + **Collaborative Innovation** + **Impact Network**



https://www.core-institute.org/

SDSC SAN DIEGO SUPERCOMPUTER CENTER

**İlkay Altıntaş, PhD** (ialtintas@ucsd.edu )

UC San Diego™
HALICIOĞLU DATA SCIENCE INSTITUTE

# 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 | 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 )

SDSC SAN DIEGO SUPERCOMPUTER CENTER

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



**İlkay Altıntaş, PhD (**ialtintas@ucsd.edu **)**

# Actionable Open Fire Science and AI:

## Right Model and Right Data
## for the Right Decision Support Workflow
## at the Right Time
## with the Right Communication

…before, during, and after a fire.

**İlkay Altıntaş, PhD** (ialtintas@ucsd.edu )

SDSC SAN DIEGO SUPERCOMPUTER CENTER

UC San Diego

# Where are we headed at WIFIRE?

- **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

- **Community Data Platforms:** WIFIRE's *Wildfire Science & Technology Commons* and *Wildfire and Landscape Resilience Data Hub* to develop standards, tools and techniques to share data and data-driven models with partners including NIST, CAL FIRE, and SDGE

- **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

- **Immersive Fire Environment:** WIFIRE's *Immersive Forest Project* leverages the AI-readiness of scientific data for new modes of teaching, training, decision-making, and public communication,

- Our platforms and products are fueled by over a dozen research projects and partnerships focused on *moving science to practice*

wifire.ucsd.edu/   SDSC SAN DIEGO SUPERCOMPUTER CENTER   UC San Diego HALICIOĞLU DATA SCIENCE INSTITUTE   ((•WIFIRE•))

# 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 real-time 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.

**REACTIVE**

## BurnPro3D

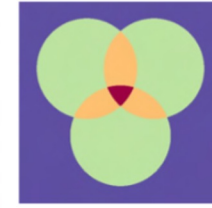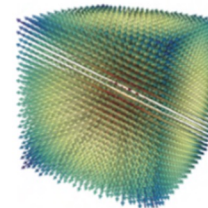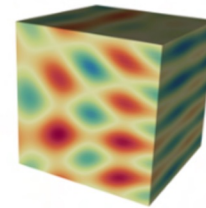In alignment with the nation's goal to increase fuel treatments to reduce wildfire risk, BurnPro3D is designed to support the preparation of burn plans as well as the implementation of prescribed burns. The interface allows burn bosses to create and visualize high-resolution 3D fire simulations and compare 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.

**PROACTIVE**

**İlkay Altıntaş, PhD** (ialtintas@ucsd.edu )

# Data and Computing Platforms

## 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.

# 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

# AI and Computing Needs for Dynamic Data-Driven Fire Modeling

-- Characterizing the dynamic fire environment : Variation of wind, smoke, moisture, fuels, fire perimeter, …

-- Detection of fire ignitions

-- Decision support for fire management

-- Prediction of potential fire ignitions

İlkay Altıntaş, PhD (ialtintas@ucsd.edu )

# AI Techniques to Condition Data and Improve Model Accuracy

**3D vs 2D**

**900x more detailed**

*Collaboration with Rod Linn (LANL), Kevin Hiers (TTRS) and Russ Parsons (USFS)*

**İlkay Altıntaş, PhD** (ialtintas@ucsd.edu )

# AI Techniques to Improve Decision Making

**Weather**

**Ignition Patterns**

**Smoke**

## PHYSICS-GUIDED MACHINE LEARNING
*To improve predictive fire behavior models*

## OPTIMIZATION
*To address complex tradeoffs and prioritization*

## EXPLAINABLE AI
*To increase scientific understanding and interpretability all along the decision-making chain*

**İlkay Altıntaş, PhD** (ialtintas@ucsd.edu )

SDSC SAN DIEGO SUPERCOMPUTER CENTER

UC San Diego ™
HALICIOĞLU DATA SCIENCE INSTITUTE

# AI in Science Communication

Visualization of multiple terrestrial LiDAR scans in the Immersive Forest prototype



Location 36.516667, -121.943056
Date: 18 June 2023
Time 8:40a
Tree Count:
Shrub Count:
Fine Fuel Volume:

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

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



İlkay Altıntaş, PhD (ialtintas@ucsd.edu )

# Immersive Fire Digital Twin

WIFIRE's **Immersive Forest** leverages the AI-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.



*Immersive Forest*
*Terrestrial & Aerial Lidar*



*Immersive Forest*
*Terrestrial Lidar*

İlkay Altıntaş, PhD (ialtintas@ucsd.edu )

wifire.ucsd.edu/

# Many Scientific Data Types in Digital Twin



**Topologies:** 1D/2D/3D - Uniform, Rectilinear, Structured, Unstructured, Polygonal, Polyhedral, AMR

**Fields:** Scalar, Vector, Multi-material

# Immersive Forest
## for
## Multimodal Communication

*Terrestrial LiDAR contextualized within Aerial scan*

# LiDAR Processing & Visualization





I. Moreno, I. Nealey, D. Roten, M. Nguyen, D. Crawl, K. O'Laughlin, M. Floca, S. Pokswinski, and İ. Altıntaş, "Visualization and Labeling of Terrestrial LiDAR Data for Three-Dimensional Fuel Classification," in Proceedings of the IEEE eScience 2023 Conference, 2023, pp. 1-2. doi: 10.1109/e-Science.58273.2023.10254841.

İlkay Altıntaş, PhD (ialtintas@ucsd.edu )

# Knowledge Representation: Spatiotemporal Data and its Challenges

Multiscale Spatial Challenge



pre-fire

post-fire

Temporal Distribution Challenge

Satellite Imagery ~30m

Aerial LiDAR ~0.5m

Terrestrial LiDAR ~1cm

İlkay Altıntaş, PhD (ialtintas@ucsd.edu )

## Knowledge Representation

## Interactive Scene Generation

### Data System

data model and underlying system for scientific data visualization

efficient query performance

high recall

empirically derived architecture

### Contextualization Engine

projections & transformations

dimensionality reduction

derivative dataset generation

add semantics and AI-driven insights

### Presentation Engine

interactive rendering

dynamic visual storytelling

data system insertion mechanism



Data ➔ Information ➔ Knowledge ➔ Interpretation and Stories

# This type of work needs the CORE4 building blocks.

**CIC⊙RE**
Cyberinfrastructure | Convergence Research | Education

## CORE4 Building Blocks

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

# AI-Integrated Applications at the Digital Continuum

# AI in Science and Research

İlkay Altıntaş, PhD (ialtintas@ucsd.edu )

# 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

**AI-Powered analysis of molecular simulations**
High-affinity generative model for target proteins

**Data-driven development of neural-network potentials from quantum chemistry data**
ML model to be used as a surrogate for expensive high-level 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

**The relationship between life span of the plant roots microscopy data and wildfire**
Deep learning model to estimate life span

**Small coronary artery calcium detectability**
Deep learning model to segment and visualize chambers of the heart

**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 iteration

**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 cell-cell communication

**İlkay Altıntaş, PhD** (ialtintas@ucsd.edu )

SDSC SAN DIEGO SUPERCOMPUTER CENTER

UC San Diego
HALICIOĞLU DATA SCIENCE INSTITUTE

# AI in Science Readiness

## "not just science + AI methods"

- Data federation and hubs
- Data quality and volume
- Knowledge management
- Benchmarks
- Scalability up and do
- Workflow managem
- Software integration and engineering
- Dev ops (also called AI ops and data ops)
- Interpretability and explainability
- Workforce training and culture/incentive building

**Requires a full team and enabling integrated data platforms**

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

**Workflow integration requires a digital continuum composed through:**

- system federation
- reusable capability services
- solutions integrating services

İlkay Altıntaş, PhD (ialtintas@ucsd.edu )

# AI in science requires data and knowledge hubs including:

- data federation
- knowledge management
- readily available standard data services
- equitable access

**İlkay Altıntaş, PhD** (ialtintas@ucsd.edu )

# 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.

İlkay Altıntaş, PhD (ialtintas@ucsd.edu )

**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*

İlkay Altıntaş, PhD (ialtintas@ucsd.edu )

SDSC SAN DIEGO SUPERCOMPUTER CENTER

UC San Diego™
HALICIOĞLU DATA SCIENCE INSTITUTE

Big Data and IoT ↔ Artificial Intelligence ↔ Modeling and Simulation

Capability ↔ Capacity

Big Data

xPU → GPU, CPU, TPU, IPU, QPU,…

Edge  FPGA

Cloud, HPC, Storage

# Some Composable Systems

İlkay Altıntaş, PhD (ialtintas@ucsd.edu )

**SDSC** SAN DIEGO SUPERCOMPUTER CENTER

# E X P A N S E

## 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

### REMOTE CI INTEGRATION

CLOUD

Open Science Grid

Heterogeneous Resources

### 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

**SDSC** SAN DIEGO SUPERCOMPUTER CENTER

**İlkay Altıntaş, PhD** (ialtintas@ucsd.edu )

**UC San Diego** HALICIOĞLU DATA SCIENCE INSTITUTE

Expanse Composable Systems Framework

# National Research Platform



https://nationalresearchplatform.org/

# First composable cluster is federated!

## EXPANSE (Enthalpy) + CHASE-CI (Nautilus)

**EXPANSE (Enthalpy)**

**DATA LIFECYCLE MANAGEMENT**
*e.g., active data repositories, long-term archives, knowledge networks, data reuse services*

Systems and services are only useful if 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



- Edge-Cloud Unified Environment for prototyping AI models to deploy on the Edge

- 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

**İlkay Altıntaş, PhD** (ialtintas@ucsd.edu )

# Fire Simulations using Composable Systems and Edge Smoke Detection

- Three workflows
  - Smoke – Sage Edge App
  - Fire simulator
  - AI Training

- Both the fire simulator and training workflows are can be run on Expanse or Nautilus through the federation layer

İlkay Altıntaş, PhD (ialtintas@ucsd.edu )

**RESPONSIBILITY**
*e.g., accuracy, privacy, explainability, ethics*

**REPRODUCIBILITY**

**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.**

**RESPONSIBILITY**
*e.g., accuracy, privacy, explainability, ethics, equity*

**REPRODUCIBILITY**

**TEAM SCIENCE**

**DATA LIFECYCLE MANAGEMENT**
*e.g., active data repositories, long-term archives, knowledge networks, data reuse services*

## 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*

**İlkay Altıntaş, PhD** (ialtintas@ucsd.edu )

SDSC SAN DIEGO SUPERCOMPUTER CENTER

UC San Diego ™
HALICIOĞLU DATA SCIENCE INSTITUTE

# Use-Inspired Composability from Systems to Services

**Team Science and Responsible Solutions**

**Services and Applications**

↕

**Federation Software and APIs**

↕

**Data and Compute Infrastructure**

**Data and Knowledge Systems**

- 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

# Democratization of CI and Data Access

İlkay Altıntaş, PhD (ialtintas@ucsd.edu )

# 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

# How do we bridge the data gaps?

http://www.nationaldataplatform.org

**İlkay Altıntaş, PhD** (ialtintas@ucsd.edu )

http://www.nationaldataplatform.org

İlkay Altıntaş, PhD (ialtintas@ucsd.edu )

# What is the National Data Platform?

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

**FOCUS AREAS:**

- 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

http://www.nationaldataplatform.org



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

# 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 Services

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**

**İlkay Altıntaş, PhD** (ialtintas@ucsd.edu )

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

**NDP Federation**



Workflow composition currently via API and/or Python client library

**NDP POP Factory**

deploys — deploys

NDP POP service stack customization

## Standard NDP Services

API + Library

- Catalog
- Search
- Data Democratization

(...)

## sciDX Stack

API + Library

- Data Staging service
- Data Streaming service
- In-situ Processing

(...)

**Third-Party Stacks**

**User-defined Services**

scidx 0.2.0
Python client library for interacting with the sciDX API
Aug 2, 2024

scidx-tools 0.1.0
Python client library for complementing the sciDX library
Jul 1, 2024

**Deployment models:** docker    kubernetes    Cloud    HPC

-- Standalone --    -- Scalable (cluster) --

# 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 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

# 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)**
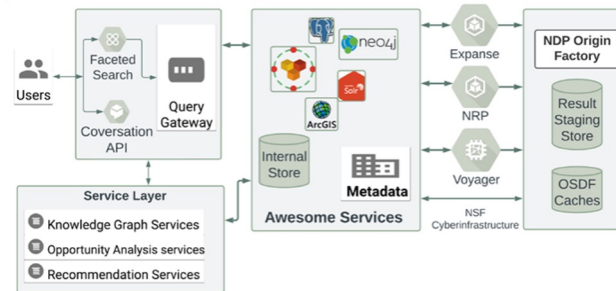
60

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

# 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.

**INTRODUCING THE**

# WILDFIRE TECHNOLOGY COMMONS

We believe that avoiding devastating wildfires requires urgent, innovative, and collaborative solutions. The Wildfire Technology Commons is a bold new initiative designed to accelerate technological innovations for wildfire management and mitigation. We are building a community platform around open data, cutting-edge science, AI, and shared knowledge.

https://www.wildfirecommons.org

**JOIN THE NETWORK**

**CONTRIBUTE DATA & MODELS**

**BECOME A PATHFINDER**

**NIST**
**National Institute of Standards and Technology**
U.S. Department of Commerce
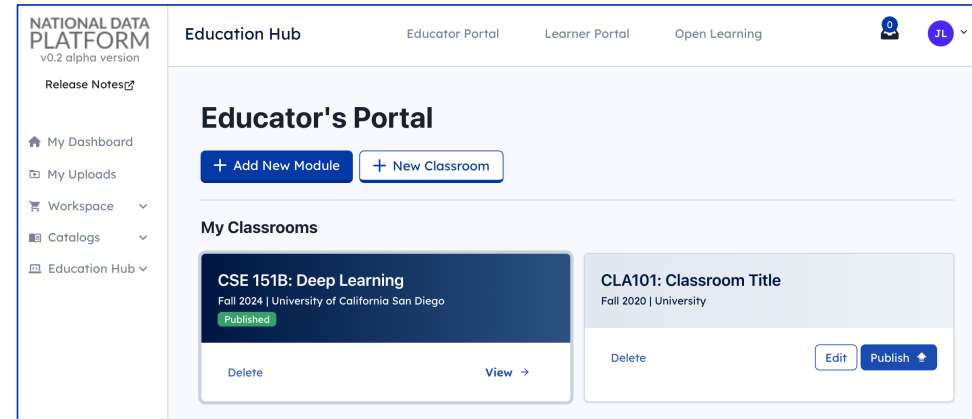
**NDP Data Challenges** for students and researchers



**lub** to
nts access
ecosystem



NATIONAL DATA PLATFORM

Designed to ensure that we are developing broadly accessible services for equitable education and community building.

The challenge questions require using data and models in an environment that requires computing and big/large data stores, which would typically be unavailable to a student or researcher without the NDP Education Hub.

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.

NATIONAL DATA PLATFORM
v0.2 alpha version

Release Notes⬀

Education Hub     Educator Portal     Learner Portal     Open Learning

🏠 My Dashboard
🖼 My Uploads
🛒 Workspace ▾
📖 Catalogs ▾
📺 Education Hub ▾

**Educator's Portal**

[ + Add New Module ]   [ + New Classroom ]

**My Classrooms**

| CSE 151B: Deep Learning | CLA101: Classroom Title |
|---|---|
| Fall 2024 \| University of California San Diego | Fall 2020 \| University |
| **Published** | |
| Delete    View → | Delete    Edit   Publish ⬆ |

# Education and capacity building through data challenges

# To sum up…

Emerging new applications require integrated AI in dynamically composed workflows, but there are significant data gaps to be addressed.

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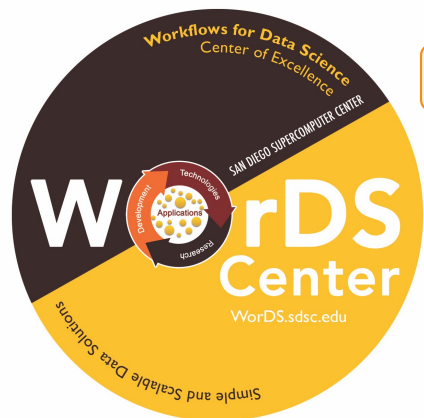


**Embrace Complexity!**

## 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

**İlkay Altıntaş, PhD** (ialtintas@ucsd.edu )

Contact:  Ilkay Altintas, Ph.D.     Email: ialtintas@ucsd.edu

https://words.sdsc.edu/

https://wifire.ucsd.edu/

**We are hiring!**
https://www.sdsc.edu/about_sdsc/careers.html

# Questions?

The presented work is collaborative work with many wonderful individuals,
and parts of it are funded by various government agencies, UC San Diego and various industry, government and foundation partners.

İlkay Altıntaş, PhD (ialtintas@ucsd.edu )