# National Data Platform (NDP) as an AI Research Resource for All

Presentation and LLM as a Service Tutorial at the 5NRP Meeting
March 19, 2024 - San Diego

**İlkay ALTINTAŞ, Ph.D.**

**University of California, San Diego**

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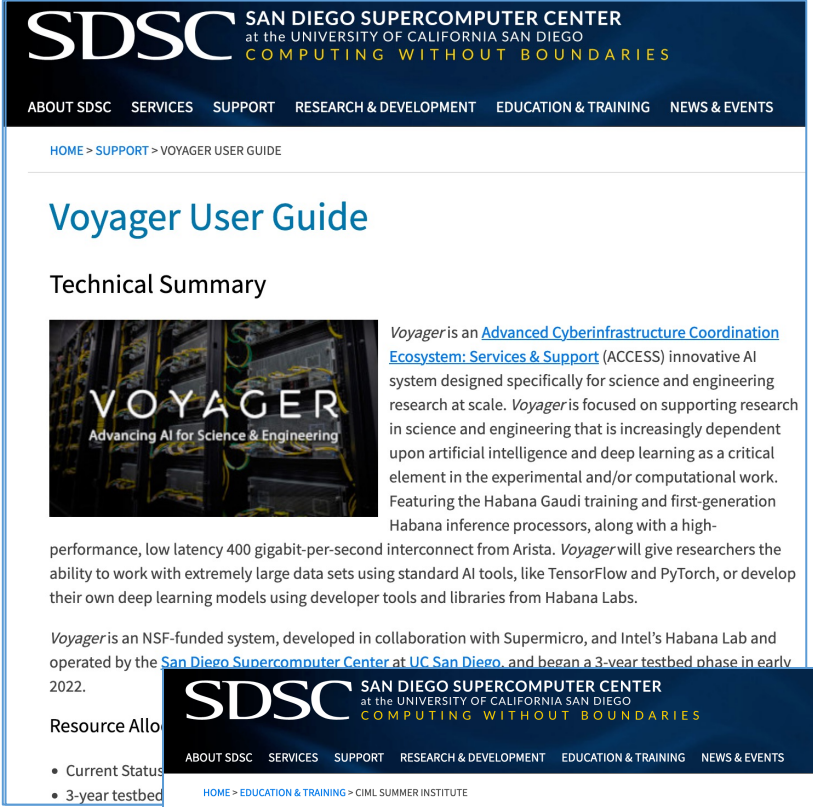Founding Fellow, **Halıcıoğlu Data Science Institute**
Founding Director, **Workflows for Data Science Center of Excellence**
Founding Director, **WIFIRE Lab**

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

SDSC SAN DIEGO SUPERCOMPUTER CENTER

UC San Diego™
HALICIOĞLU DATA SCIENCE INSTITUTE

# Some Examples of AI at SDSC

- Exploring new architectures in support of AI in research and engineering
- Teaching best practices for machine learning and data science applications
- Building methods for AI-integrated societal impat
- NAIRR Pilot Activities

---

**SDSC SAN DIEGO SUPERCOMPUTER CENTER**
at the UNIVERSITY OF CALIFORNIA SAN DIEGO
COMPUTING WITHOUT BOUNDARIES

ABOUT SDSC    SERVICES    SUPPORT    RESEARCH & DEVELOPMENT    EDUCATION & TRAINING    NEWS & EVENTS

HOME > SUPPORT > VOYAGER USER GUIDE

## Voyager User Guide

### Technical Summary

*Voyager* is an Advanced Cyberinfrastructure Coordination Ecosystem: Services & Support (ACCESS) innovative AI system designed specifically for science and engineering research at scale. *Voyager* is focused on supporting research in science and engineering that is increasingly dependent upon artificial intelligence and deep learning as a critical element in the experimental and/or computational work. Featuring the Habana Gaudi training and first-generation Habana inference processors, along with a high-performance, low latency 400 gigabit-per-second interconnect from Arista. *Voyager* will give researchers the ability to work with extremely large data sets using standard AI tools, like TensorFlow and PyTorch, or develop their own deep learning models using developer tools and libraries from Habana Labs.

*Voyager* is an NSF-funded system, developed in collaboration with Supermicro, and Intel's Habana Lab and operated by the San Diego Supercomputer Center at UC San Diego, and began a 3-year testbed phase in early 2022.

### Resource Alloc

- Current Status
- 3-year testbed

---

**SDSC SAN DIEGO SUPERCOMPUTER CENTER**
at the UNIVERSITY OF CALIFORNIA SAN DIEGO
COMPUTING WITHOUT BOUNDARIES

ABOUT SDSC    SERVICES    SUPPORT    RESEARCH & DEVELOPMENT    EDUCATION & TRAINING    NEWS & EVENTS

HOME > EDUCATION & TRAINING > CIML SUMMER INSTITUTE

## CIML Summer Institute

**Applications for the CIML Summer Institute 2024 is now open!
Apply today!**

**Application deadline: Friday, April 12, 2024**

- **Preparation Day (virtual):** Tuesday, June 18, 2024
- **Summer Institute (in-person):** Tuesday, June 25 – Thursday, June 27, 2024
- Location: SDSC Auditorium, UC San Diego

The San Diego Supercomputer Center (SDSC) **Cyberinfrastructure-Enabled Machine Learning (CIML)** project is focused on teaching researchers and students the best practices for effectively running machine learning (ML) and data science applications on advanced cyberinfrastructure (CI) and high-performance computing (HPC) systems.

The CIML Summer Institute introduces machine learning (ML) concepts to researchers, developers and educators to techniques and methods needed to migrate their ML applications from smaller, locally run resources, such as laptops and workstations, to large-scale HPC systems, such as the SDSC's Expanse supercomputer. Participants will have the opportunity to accelerate their learning process through highly

---

**CORE INSTITUTE**

## Data and AI Tools for Regional Food Systems

### 2024 CORE FELLOWS CALL FOR APPLICATIONS

The CORE Institute is excited to announce a call for applications for the 2024 CORE Fellows Program.

- When? May **16-17, 2024**
- Where? **UC San Diego**
- We welcome applicants with experience and/or interest in regenerative agriculture, food security, data-driven technologies for sustainable supply chains, or related fields.

The Fellowship includes funding for travel expenses. Additionally, Fellows will have the option to submit workshop proposals for the CORE Summit (Fall 2024). Selected proposals will be funded by the CORE Institute.

- **Apply online by March 21st, 2024.**

Visit the CORE Institute website to subscribe, apply, or learn more.

core-institute.org

**SDSC@UC San Diego** The CORE Institute is guided by a Leadership Council composed of members who have participated in the NSF's Convergence Accelerator. The program is managed by the San Diego Supercomputer Center at UC San Diego.

---

# Equity of AI in Research

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

FROM:      Dr. Alondra Nelson
            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 and data.

**The case for open data**

**Empowering citizens & strengthening accountability**
- *Promotes more accountability*
- *Increases citizen engagement*

**Innovation & efficiency in government agencies**
- *Decreased workloads*
- *Inter-agency collaboration*
- *Improved policy design*

**Creating wider value for the economy**
- *Open data creates value added services for the entire economy*

*OECD*

# The Minds We Need

**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**

Yubo Qin, *Rutgers University, New Brunswick, NJ, 08901, USA*
Ivan Rodero and Manish Parashar, *University of Utah, Salt Lake City, UT, 84112, USA*

*Data collected by large-scale instruments, observatories, and sensor networks (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 and discovery.*

**The Missing** *Millions*

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

**October 3, 2021**
NSF OAC 2127459

# Democratization of CI and Data Access

# Open Questions for Equitable Open Research

What are the foundational data abstractions, catalogs, multipurpose services and expandable workflows for data-driven and AI-integrated application patterns?

**How can everyone effectively access and utilize these abstractions and services**?

How can services and workflows be developed and deployed on top of production-ready CI?

**How can equity be ensured for all to access & use CI** from storage to the edge-to-HPC computing continuum?

What are the governance and open science, open data and open CI requirements and challenges?

What are the required guardrails for protecting privacy, civil rights and civil liberties that will **ensure a more equitable use of such data systems and services for everything from education to new AI training and application development**?

**FOUNDATIONAL ABSTRACTIONS, CATALOGS, AND SERVICES**

**EQUITABLE OPEN CI USE**

**NEEDS, REQUIREMENTS AND CHALLENGES**

# Architecting for Equity of Research Workflows for All

- Involve diverse users in architecting around access, use, expertise and education gaps
- Improve the experience of working with data
  - e.g., serve data and knowledge systems around it
- Create an ecosystem approach to capacity building
  - e.g., through services, platforms, education of many types
- Incubate use-inspired solutions to scale
- Explore new models of allocation
  - e.g.,  service unites, credits, tokens, aggregated workflow coops
- Develop models of sustainability and scale
  - e.g., public/private partnerships, NGOs, consortiums, cooperatives

# NATIONAL DATA PLATFORM

**Addressing the Missing Middle for AI-enabled Data-driven Research and Education Workflows**

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

# National Data Platform Pilot: Services for Equitable Open Access to Data

nationaldataplatform.org

National Data Platform is a federated and extensible data and service ecosystem to promote collaboration, innovation and equitable use of data on top of existing cyberinfrastructure capabilities.

NDP enables AI-integrated science workflows that foster discovery, decision-making, policy formation and societal impact related to wildfire, climate, earthquake and food security among others.

Link to the award abstract: https://www.nsf.gov/awardsearch/showAward?AWD_ID=2333609

**NDP Hub**

Discover and use interconnected data hubs and user-facing services deployed on CI

**NATIONAL DATA PLATFORM**

**NDP Platform**

Develop and deploy services, application workflows and educational challenges

# Fostering scientific understanding, decision-making, policy formation and societal impact

**Focus:** use of existing data repositories to scientists and nonexpert users, making technology accessible to those without access to data and AI expertise

**Objectives:**
- contribute to a more equitable data and AI research
- build a broadly accessible data ecosystem
- enable diversity in
  - data sources
  - perspectives and experiences of students and researchers
  - research practices and governance processes
- equitably manage both the benefits and risks related to AI

**Case studies** initially focused on earth sciences and food security but designed to be generalizable.

NATIONAL DATA PLATFORM

**Create reusable capabilities and amplify the value of existing data repositories to benefit science, society and education.**

# Reference Architecture

Links data and cyberinfrastructure with domain specific platforms to enable value add services and open educational capabilities.



National Data Platform (NDP) landing page

# Example NDP NAIRR Integration



**NDP Hub**
- Find NAIRR resources
- Learn from notebook examples
- Use NAIRR allocations

**NDP Platform**
- Build new models and data
- Publish via GitHub and HuggingFace
- Create narrative stories and educational resources for others

# 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.

# Community expansion and stakeholder engagement

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



**Reports** — Year 3 2nd Half
- Reference Architecture
- Co-Design Workshop Outcomes
- Needs Assessment Final Outcomes

Workshop 4 — Year 3 1st Half

**Expansion**
- Co-Design with External Community
- Needs Assessment Review

**Open Knowledge Networks**
- Co-Design with NOURISH Community
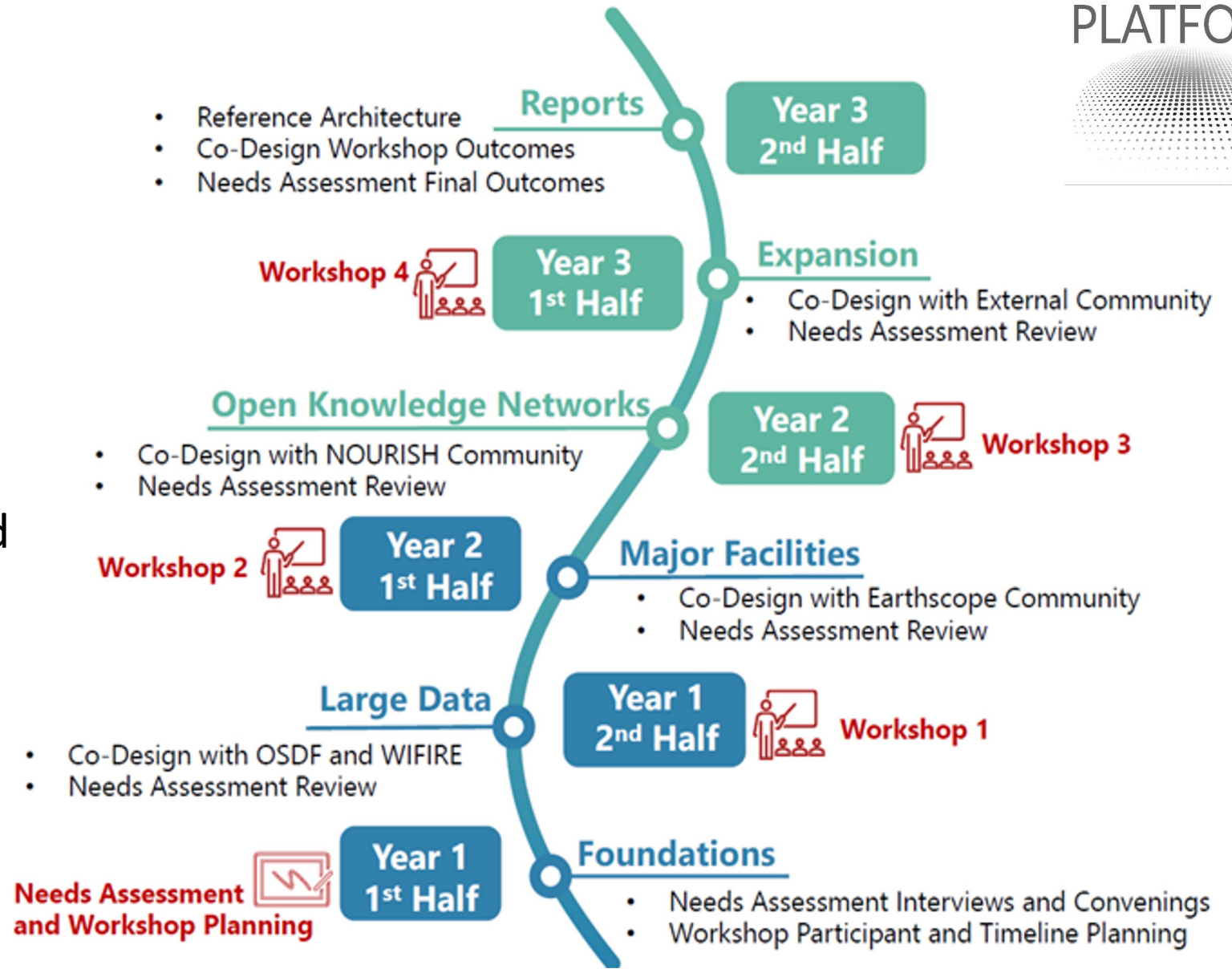- Needs Assessment Review

Year 2 2nd Half — Workshop 3

Workshop 2 — Year 2 1st Half

**Major Facilities**
- Co-Design with Earthscope Community
- Needs Assessment Review

**Large Data**
- Co-Design with OSDF and WIFIRE
- Needs Assessment Review

Year 1 2nd Half — Workshop 1

Needs Assessment and Workshop Planning — Year 1 1st Half

**Foundations**
- Needs Assessment Interviews and Convenings
- Workshop Participant and Timeline Planning

# Education and capacity building through data challenges

**NDP Data Challenges** for students and researchers

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

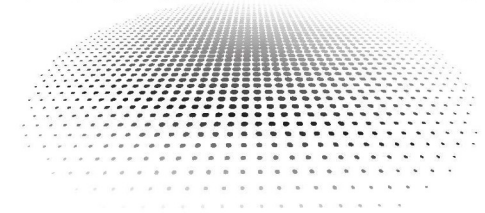**NDP Education Gateway** to provide participants access to the NDP data ecosystem

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.

NATIONAL DATA PLATFORM

**Three Co-Design Workshops**

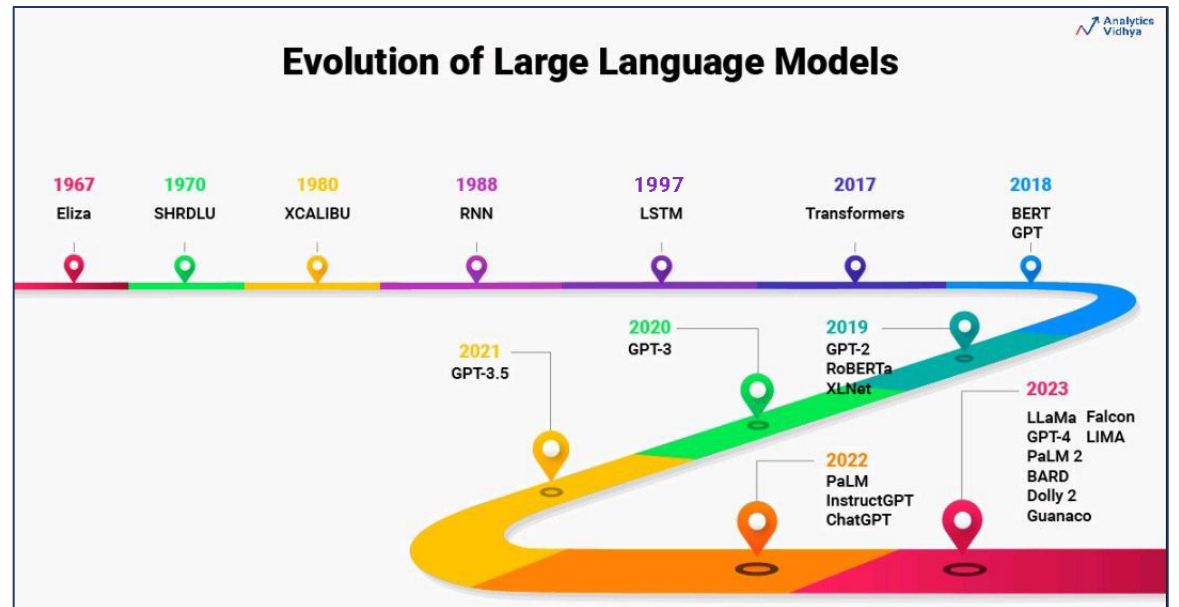Each will include a breakout session to develop a data challenge question specific to large data (W1); streaming data (W2); and graph data (W3).

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.

# An NDP Service Example:
## Generative AI and Large Language Models (LLM)

- Huge generative potential

- Ability to create human-like outputs

- Integration with complex models

- Libraries advanced technologies

  - e.g., GPT, Prompt Engineering, and vector storage

- Shortcomings on domain expertise

- Need domain-specific LLMs

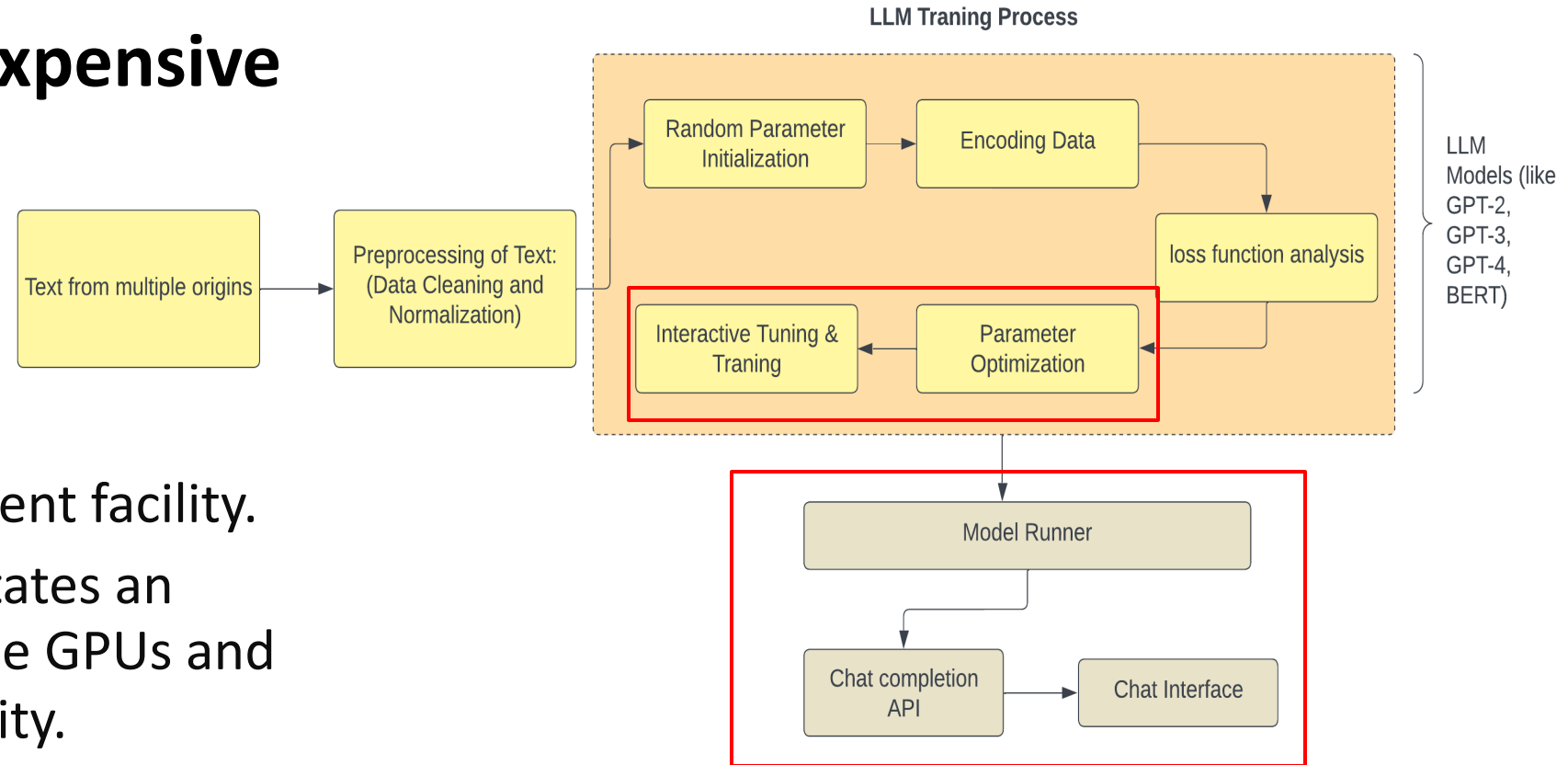  - with human-curated data and controlled knowledge



**Source:** https://www.analyticsvidhya.com/blog/2023/07/beginners-guide-to-build-large-language-models-from-scratch

"Generative AI helped workers avoid awful ideas, but it also led to more average ideas"
         - **Harvard Business Report (March – April 2024)**

# Accessing and Using LLMs is an Equity Issue

## LLM Deployment is Expensive

- Even tuning an LLM can incur substantial costs, necessitating 4-5 AT100 GPUs, expansive nodes, and an equipped deployment facility.

- Operating an LLM necessitates an infrastructure with multiple GPUs and substantial memory capacity.

**LLM Traning Process**



http://www.nationaldataplatform.org

# NDP
# LLM as a Service



- Tailored Model Selection
- Enhanced Data Control
- Privacy and Security
- Cost Efficiency
- OpenAI API and LangChain Support

## LLM Client Service

- Use an existing model
- Add context with domain-specific documents

## LLM Training Service

- Fine-tune an existing model to create a new model
- Use a larger corpus for training
- Deploy as a service

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

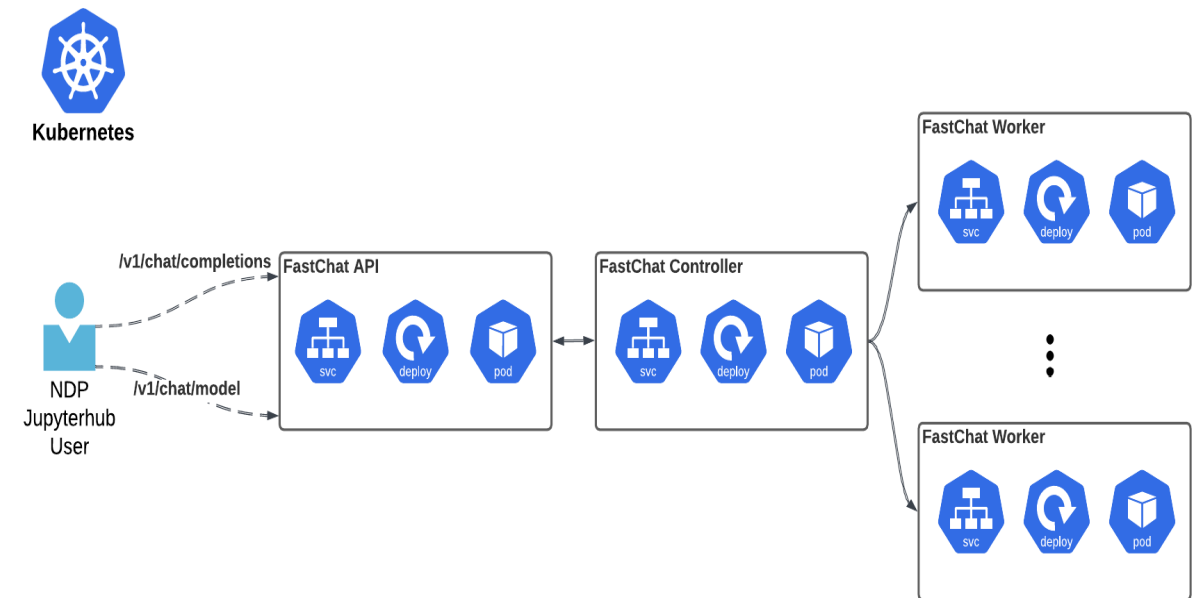# NDP LLM as a Service

## Alignment with NAIRR Objectives

o **Capacity** to support many users with a spectrum of backgrounds

o **Capabilities**

- Ability to train *(and use)* resource-intensive AI models on CI resources
- Ability to make use of a mix of computational resources
- Option to select which resources to use through a range of mechanisms, including ... optionally interactive "notebook"-like environments
- A NAIRR system should include at least one large-scale machine-learning supercomputer capable of training 1 trillion-parameter models

*In today's tutorial we are using a model with*
***7B** parameters running on NRP*

# NDP LLM Deployment Architecture

- FastChat
  - Open source LLM execution library
  - Deployed on Nautilus
    - API Server
    - Controller
    - Worker (serves different or the same LLMs)
- Currently all workers are serving the following LLMs
  - eci-io/climategpt-7b,
  - ECarbenia/grimoiresigils
  - text-embedding-ada-002

# Part 2 Comin Up at 4:30pm:

# NDP LLM-as-a-Service on NRP Tutorial

# How can we work with you?

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