**National Institute of Diabetes and Digestive and Kidney Diseases**

*Central Repository*

NIDDK-CR Resources for Research

# Data Science and Meet the Expert Webinar Series

January 21, 2026

# NIDDK Central Repository Overview

National Institute of Diabetes and Digestive and Kidney Diseases
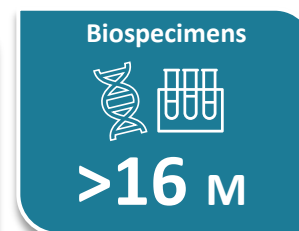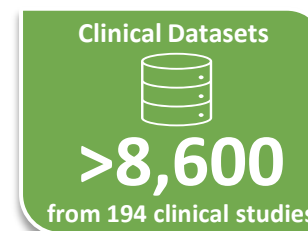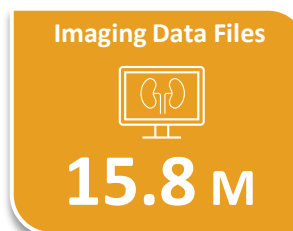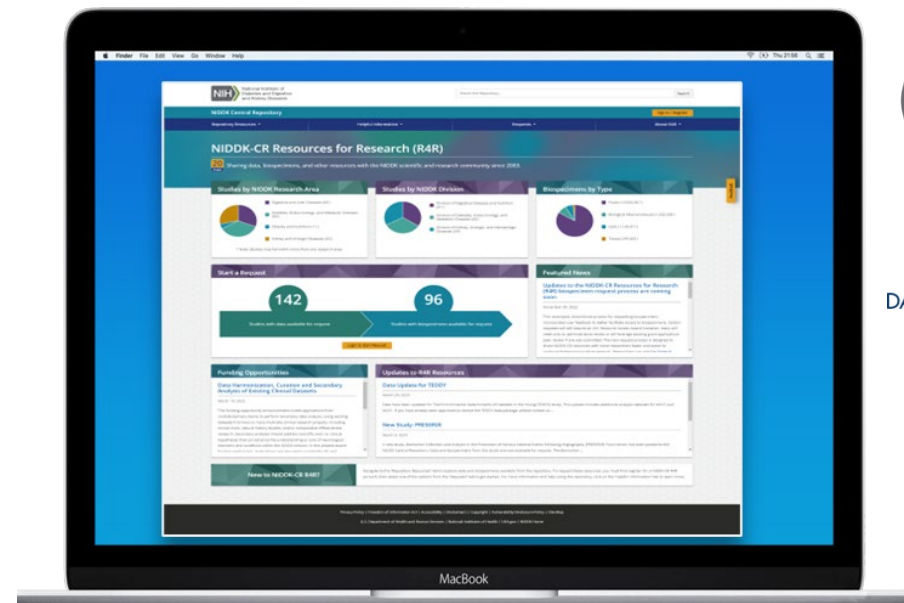
*Central Repository*

## Mission

Established in 2003 to **facilitate sharing of data, biospecimens, and other resources** generated from studies supported by NIDDK and within NIDDK's mission by making these **resources available for request to the broader scientific and research community**.

- Supports receipt and distribution of data and biospecimens in a manner that is ethical, equitable, and efficient

- Enables investigators not involved with the original work to test new hypotheses without the need to collect new data or biospecimens

- Promotes FAIR (Findable, Accessible, Interoperable, and Reusable) and TRUST (Transparency, Responsibility, User focus, Sustainability, and Technology) principles

**Recorded past tutorials, webinars, and other educational resources can be found on the NIDDK-CR website**

| | | |
|---|---|---|
| **Imaging Data Files** | **Clinical Datasets** | **Biospecimens** |
| 15.8 M | >8,600 from 194 clinical studies | >16 M |
| **Registered Users** | **Weekly Users** | **Public Releases** |
| 7,026 | >5,000 | >875 |

# Analytics Workbench Functionality

**National Institute of Diabetes and Digestive and Kidney Diseases**

*Central Repository*

***Streamlining*** *end-to-end data science lifecycle and discovery of data-driven biomedical insights.*

## Innovation and ease of use

A cloud-based analytics environment where researchers and data scientists can access a suite of integrated analytics tools and cloud computing resources to participate in data challenges and AI innovation.

### Expected Benefits of Analytics Workbench:

| | | |
|---|---|---|
| **Promote Collaboration** | **Support AI Innovation** | **Minimize Data Movement** |
| **Improve User Experience** | **Discover Data Insights** | **Advance NIDDK Research Mission** |

---

**Insights and Interactions**

| Search | Access Services | Communication Services | Collaborate & Share | Chatbot | Data Insights |
|---|---|---|---|---|---|

**Analytics Workbench**

**Data Analytics and AI Services**

**Data Discovery**
- Search
- Query
- Data Discovery and Exploration

**Descriptive**
- Usage Reports
- Dashboards
- Self Service BI

**Predictive**
- Predictive modeling with biomarkers
- Disease Forecasting
- Simulation

**Cognitive AI**
- Natural Language Processing
- Machine Learning
- Generative AI

**User-Centric**
- Personalization
- Collaboration workspace
- Communication

**Prescriptive Data Analytics and AI Tools**

jupyter | SAS | GitHub Copilot

**Data Services / API Connectors**

**Data and Cloud Infrastructure Foundation**

Internal Data Sources | External Data Sources

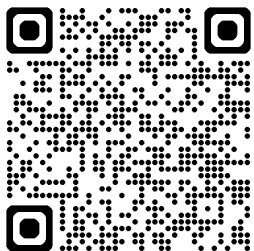**Platform Provisioning | Deployment | Monitoring | Security | Access | User Management**

**National Institute of Diabetes and Digestive and Kidney Diseases**
*Central Repository*

## Goals of NIDDK-CR Data-science centric challenge series

- Develop tools, approaches, models and/or methods to increase data interoperability and usability for artificial intelligence (AI) and machine learning (ML) applications

- Augment and enhance existing data for future secondary research, including data-driven discovery by AI/ML researchers

- Discover innovative approaches to enhance the utility of datasets for AI/ML applications

**Visit our website for more information on our data-centric movement and to learn more about our past data-challenges**

# Meet the Experts

Chen Li is a professor in the Department of Computer Science at UC Irvine. He received his Ph.D. degree in Computer Science from Stanford University, and his M.S. and B.S. in Computer Science from Tsinghua University, China, respectively. He was a recipient of an NSF CAREER award and several test-of-time publication awards, a part-time visiting research scientist at Google, an ACM distinguished member, and an IEEE fellow. He was a co-founder and CTO of a startup to commercialize his research.

Kun Woo Park and Jiadong Bai are PhD students in Computer Science at UC Irvine working on the Texera project.

# Advancing Collaborative Data Science with Texera

Prof. Chen Li

with Kun Woo Park and Jiadong Bai

Department of Computer Science, UC Irvine

January 21, 2026

# Apache Texera: Overview

- Supporting data science and AI/ML as workflows

- Cloud services (no installation, software patches)

- Supporting community-based sharing of data and workflows

- Shared editing/execution

- Supporting Python, R, Java as user-defined functions (UDFs)

- Started in 2016

- Open source (being incubated by Apache)

- Parallel engine, scalable

Chen Li, UCI

# Part of NIDDK dkNET Computational Core

Making Data Science and AI/ML easily available to the NIDDK community

# Example application: sequence analysis in biology



Alice: Biologist (PI)

Sally: Bioinformatician

Bob: Bioinformatician

Chen Li, UCI

**Sequence analysis pipeline**

FASTQ Sequences

Cellranger — Step 1

Count Matrix

Quality Control — Step 2

Normalization — Step 3

Dimensionality Reduction — Step 4

Data Integration — Step 5

Clustering — Step 6

Cluster Annotation — Step 7

AI/ML analysis

Trajectory analysis

Differential "gene" expression analysis

# Coding challenges

- Coding is hard!
- Version control of libraries
- Needs servers
- Slow on large data
- Not every lab can afford a bioinformatician



Data preparation
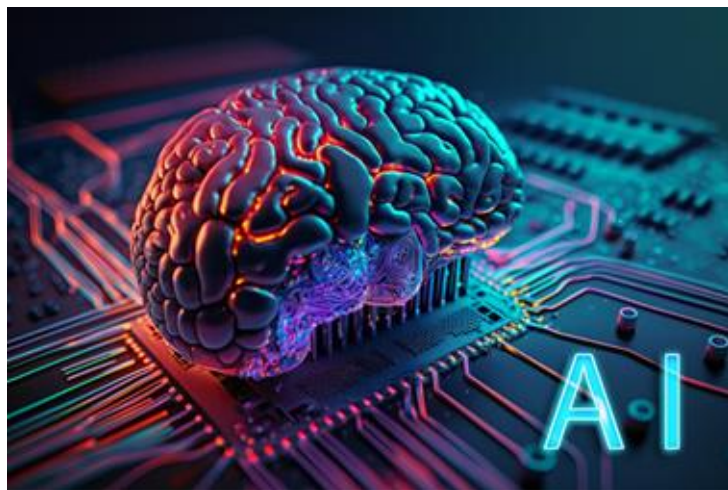
Data analytics

Visualization

Sally: Bioinformatician

# Collaboration challenges

- Collaborators of different backgrounds
  - Biologists
  - Bioinformaticians
  - Computer scientists
- Collaborators from different organizations
  - Same lab: senior students vs new students
  - Other labs

Chen Li, UCI

# AI/ML opportunities

- How to utilize state-of-the-art AI/ML technologies?

- Require advanced coding skills

- Not easily available

Chen Li, UCI

# Our solution



**Collaborative data science and AI/ML using GUI-based workflows**

Chen Li, UCI

# Open source (Apache Incubating)

# Demo!

Chen Li, UCI

**Statistics**

| ASF and GitHub (as of Dec. 2025) | |
| --- | --- |
| PPMC members & committers | 14 |
| ASF mentors | 4 |
| GitHub contributors | 148 |
| Open issues | 129 |
| Closed issues | 783 |
| Open pull requests | 32 |
| Closed pull requests | 3,101 |

| Usage and deployments (as of Dec. 2025) | |
| --- | --- |
| Users | > 600 |
| Workflows created | > 3,000 |
| Workflow versions edited | > 273,000 |
| Workflow executions | > 51,000 |
| Largest deployment: node # | 100 |
| Largest deployment: core # | 400 |

UCI

# Example: analyzing brain images, 256GB

# Teaching non-STEM students AI/ML using Texera

# 2025 dkNET Summer Bootcamp

| | | | | | |
|---|---|---|---|---|---|
| Monday, July 21, 2025 9 am-1 pm PDT | Introduction to data science, data modeling, and data preparation | Dr. Chen Li and Sarah Asad | | Getting familiar with concepts related to data science. Students will start a capstone project using provided data. | https://hub.texera.io |
| Tuesday, July 22, 2025 9 am-1 pm PDT | Python programming to do data science | Dr. Chen Li and Sarah Asad | | Using Python to do data science. Students will continue working the capstone project. | https://hub.texera.io |
| Tuesday, July 23, 2025 9 am-1 pm PDT | Introduction to machine learning | Dr. Wei Wang and Alexander Taylor | | Getting familiar with concepts related to AI/ML Students will finish the capstone project. | https://hub.texera.io |
| Friday, July 25, 2025 10 am-12 pm PDT | Discussion Session: FAIR Data and DMSP | Dr. Maryann Martone Dr. Jeffrey Grethe | Assignment discussion: 1) Based on your research project and the data it is using, use dkNET tools to help you select appropriate repository(s) and work through what is needed to manage and share your data in compliance with NIH's new DMSP requirements. Work with some of your data to ensure it is FAIR and Frictionless (https://frictionlessdata.io) - document what would be needed as part of your data collection and management practices. | 1) Check-in project progress 2) Assignment discussion (FAIR data; Data Management) | |

https://dknet.org/about/summer_course_2025

# Ongoing efforts

- Support management of ML models

- Incorporate more AI techniques to the platform

- Make analysis pipelines to the community

- Improve security and privacy

- High performance and scalability

- Elastic computing using cloud resources

- …

# Summary: Apache Texera

- Cloud-computing platform

- GUI-based workflows (no coding needed)

- Collaboration and sharing of data/analyses

- Parallel computing: for big data

- Supporting multiple languages: Python, R, Java, …

- Supporting AI/ML (training, inference, …)

Chen Li, UCI

# **Advancing Collaborative Data Science with Texera**

Prof. Chen Li

Department of Computer Science

UC Irvine