# Computational Data Science Strategies

**Getting Ready for a Data Science 101 Course**

**Deborah Duran**, PhD · NIMHD
**Luca Calzoni**, MD MS PhD Cand. · NIMHD
**Kenneth Wilkins**, PhD · NIDDK

January 17, 2024

# Thank you

**National Institute on Minority Health and Health Disparities** + **NIH Office of Data Science Strategy** + **NIH National Institute of Nursing Research**

## NIMHD

Dr. Eliseo Perez-Stable

### ODSS

Dr. Susan Gregurick

### NIH/OD

Dr. Larry Tabak

### NINR

Dr. Shannon Zenk

## NINR

Rebecca Hawes
Micheal Steele
John Grason

### ORWH

### OMH

## NIMHD OCPL

Kelli Carrington
Thoko Kachipande
Corinne Baker

### BioTeam

### STRIDES

### Terra

### SIDEM

### RLA

### Broad Institute

## CDE Working Group

Deborah Duran
Luca Calzoni
Rebecca Hawes
Micheal Steele
Kelvin Choi
Paula Strassle
Deborah Linares
Crystal Barksdale
Gneisha Dinwiddie
Jennifer Alvidrez
Matthew McAuliffe
Carolina Mendoza-Puccini
Simrann Sidhu
Tu Le

# Experience poll

**Please check your level of experience with the following:**

|  | None | Some | Proficient | Expert |
|---|---|---|---|---|
| Python | ☐ | ☐ | ☐ | ☐ |
| R | ☐ | ☐ | ☐ | ☐ |
| Cloud computing | ☐ | ☐ | ☐ | ☐ |
| Terra | ☐ | ☐ | ☐ | ☐ |
| Health disparities research | ☐ | ☐ | ☐ | ☐ |
| Health outcomes research | ☐ | ☐ | ☐ | ☐ |
| Algorithmic bias mitigation | ☐ | ☐ | ☐ | ☐ |

# ScHARe

**ScHARe is a cloud-based population science data platform** designed to accelerate research in health disparities, health and healthcare delivery outcomes, and artificial intelligence (AI) bias mitigation strategies

**ScHARe aims to fill three critical gaps:**

- Increase participation of **women & underrepresented populations with health disparities** in data science through data science skills training, cross-discipline mentoring, and multi-career level collaborating on research

- Leverage population science, SDoH, and behavioral Big Data and cloud computing tools to foster a **paradigm shift** in healthy disparity, and health and healthcare delivery outcomes research

- **Advance AI bias mitigation and ethical inquiry** by developing innovative strategies and securing diverse perspectives
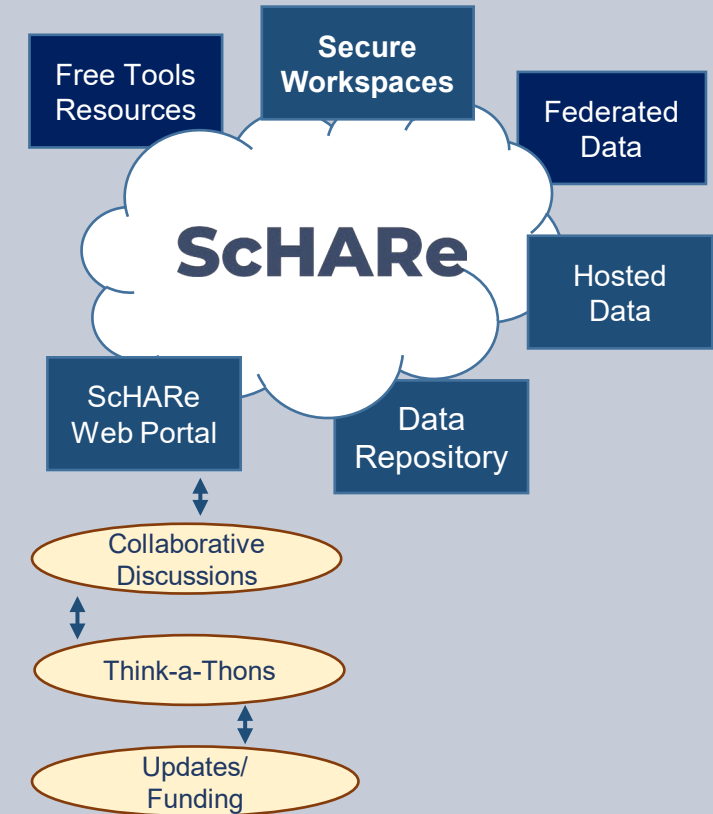


nimhd.nih.gov/schare

# ScHARe Components

ScHARe co-localizes within the cloud:

- **Datasets** (including social determinants of health and social science data) relevant to minority health, health disparities, and health care outcomes research

- **Data repository** to comply with the required hosting, managing, and sharing of data from NIMHD- and NINR-funded research programs

- **Computational capabilities and secure, collaborative workspaces** for students and all career level researchers

- **Tools for collaboratively evaluating and mitigating biases** associated with datasets and algorithms utilized to inform healthcare and policy decisions

**Frameworks**:    Google Platform, Terra, GitHub, NIMHD Web ScHARe Portal



**Intramural & Extramural Resource**

Free Tools Resources

Secure Workspaces

Federated Data

ScHARe

Hosted Data

ScHARe Web Portal

Data Repository

Collaborative Discussions

Think-a-Thons

Updates/ Funding

nimhd.nih.gov/schare

# ScHARe Data Ecosystem

Researchers can access, link, analyze, and export **a wealth of datasets** within and across platforms relevant to research about health disparities, health care outcomes and bias mitigation, including:

- **Google Cloud Public Datasets:** publicly accessible, federated, de-identified datasets hosted by Google through the Google Cloud Public Dataset Program

  **Example**: *American Community Survey (ACS)*

- **ScHARe Hosted Public Datasets:** publicly accessible, de-identified datasets hosted by ScHARe

  **Example**: *Behavioral Risk Factor Surveillance System (BRFSS)*

- **Funded Datasets on ScHARe:** publicly accessible and controlled-access, funded program/project datasets using Core Common Data Elements shared by NIH grantees and intramural investigators to comply with the NIH Data Sharing Policy

  **Examples**: *Jackson Heart Study (JHS); Extramural Grant Data; Intramural Project Data*

Datasets are categorized by content based on the CDC **Social Determinants of Health categories**:

1. Economic Stability
2. Education Access and Quality
3. Health Care Access and Quality
4. Neighborhood and Built Environment
5. Social and Community Context

with the addition of:
- **Health Behaviors**
- **Diseases and Conditions**

Users will be able to **map and link** across datasets

# ScHARe Data Ecosystem Structure

**FEDERATED PUBLIC DATA 240+**

Hosted by Google & ScHARe

**REPOSITORY**

**CDE FOCUSED**

CDEs enhances Data Interoperability (Aggregation) by using semantic standards and concept codes

*Innovative Approach:*

*CDE Concept Codes Uniform Resource Identifier (URI)*

## What is a CDE?

A common data element (CDE) is a standardized, precisely defined question that is paired with a set of specific allowable responses, that is then used systematically across different sites, studies, or clinical trials to ensure consistent data collection

# ScHARe CDEs Labels

- Age
- Birthplace
- Zip Code
- Race and Ethnicity
- Sex
- Gender
- Sexual Orientation
- Marital Status
- Education
- Annual Household Income
- Household Size

- English Proficiency
- Disabilities
- Health Insurance
- Employment Status
- Usual Place of Health Care
- Financial Security / Social Needs
- Self Reported Health
- Health Conditions (Associated Medications/Treatments)

**NIMHD Framework
**Health Disparity Outcomes

**NIH Endorsed**

(** project level CDE)

**NIH CDE Repository: https://cde.nlm.nih.gov/home**

Cross-walked with PhenX SDoH

NIH-endorsed CDEs have been reviewed and approved by an expert panel, and meet established criteria. They are designated with a gold ribbon. 🎗

# ScHARe
## REPOSITORY

**COMMON DATA ELEMENTS**

**DATA UPLOAD**

**DATA MAPPING, DOWNLOAD AND EXPORT**

### NLM CDE Repository
**Coded NIMHD Common Data Elements**

- Labels
- Questions
- Permissible Values

**A T O**

### Common Data Elements + Data

### Data Access
**Based On PII Levels and User Needs:**
- Public
- Data Use Agreement
- Private

### Acquired Google and ScHARe Hosted Datasets

| Overview |
| Data Dictionaries |
| Data Updates |

### Project and Key Acquired Datasets

**Overview**

Description and Links to Overview Material

4-Privacy Levels

**COMMON DATA ELEMENTS**

**Data**

**Metadata**

Data Dictionaries

**Analysis Ready**

### RAS Single Sign-on

### Other Cloud Platforms
**AnVil, BDC, All of Us**

### DATA MAPPING
**ACROSS DATASETS AND PLATFORMS BASED ON CDES**

EXAMPLE: CDE linked
ACS      NIMHD Project      BioData Catalyst

**Aggregated Data Set**

**CDE Linked Project Data**

**Data Download in a Variety of Formats**
CSV, TSV, XLSX

**Data Export to Terra for Analysis**
**Workspaces**

**Visualizations Tools**
**Shiny**

# ScHARe

## Project & federated dataset mapping

| |
|---|
| Project Title |
| Project Description |
| Core Common Data Elements |
| Other Project Data |
| Data Dictionary |

**+** AMERICAN COMMUNITY SURVEY

**+**

**+** Medical Expenditure Survey (MEPS)

**+**

**+** Pharmacy and health insurance databases

## Mapping across cloud platforms



ScHARe

All of Us RESEARCH PROGRAM

Terra

AnVIL

BioData CATALYST

**UPCOMING**

# ScHARe

**Repository
CDE Focused
for Data
Interoperability**

Coming Soon

# Secure workspace



- Secure workspace **for self or collaborative research**

- **Assign roles**: review or admin

- **Host own data and code**

# Notebooks analytics



# Workflows - Modular codes

- **Copy and paste analytics**



- Modular codes developed for reuse
- **Adding SAS**

# ScHARe Registrations



1900+ unique users

# Think-a-Thon Tutorials

| | |
|---|---|
| February | **Artificial Intelligence and Cloud Computing 101** |
| March | **ScHARe 1 – Accounts and Workspaces** |
| April | **ScHARe 2 – Terra Datasets** |
| May | **ScHARe 3 – Terra Google-hosted Datasets** |
| | *ScHARe for Educators (Community Colleges & Low Resource MSIs)* |
| June | **ScHARe 4 – Terra ScHARe-hosted Datasets** |
| July | **An Introduction to Python for Data Science – Part 1** |
| August | **An Introduction to Python for Data Science – Part 2** |
| | *ScHARe for American Indian / Alaska Native Researchers* |
| September | **ScHARe 5: A Review of the ScHARe Platform and Data Ecosystem** |
| October | **Preparing for AI 1: Common Data Elements and Data Aggregation** |
| November | **Preparing for AI 2: An Introduction to FAIR Data and AI-ready Datasets** |
| January | **Preparing for AI 3: Computational Data Science Strategies 101** |
| | *ScHARe for Coders and Programmers to conduct Research (Jan 31)* |

**bit.ly/think-a-thons**

**Upcoming**

**ScHARe**

# Think-a-Thons (TaT)

## Research Teams

**Title:** Data Science Projects 1 – Health Disparities and Individual SDoH

**Description:** Exploring the impact of individual Social Determinants of Health on health outcomes: a hands-on session for researchers and students at all levels interested in collaborating on ScHARe to develop innovative research questions and projects leading to publications.

**Title:** Data Science Projects 2 - Health Disparities and Structural SDoH

**Description:** Assessing the impact of structural Social Determinants of Health on health outcomes: a hands-on session for researchers and students at all levels interested in collaborating on ScHARe to develop innovative research questions and projects leading to publications.

**Title:** Data Science Projects 3 – Health Outcomes

**Description:** Investigating the influence of non-clinical factors on disparities in health care delivery: a hands-on session for researchers and students at all levels interested in collaborating on ScHARe to develop innovative research questions and projects leading to publications.
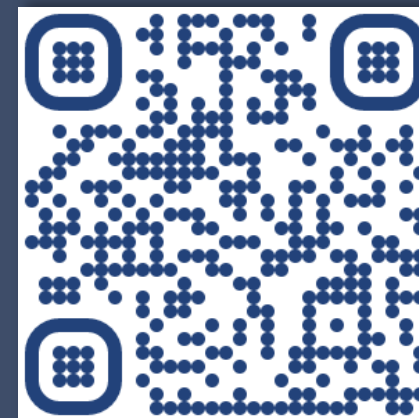
- Multi-career (students to sr. investigators)
- Multi-discipline (data scientist & researchers)
- Feature Datasets with Guest Expert Leads
- Secure experts in topic area, analytics, data sources etc. to provide guidance
- Generate research idea - decide potential design, datasets & analytics
- Select co-leads to coordinate completion outside of TaT
- Publications
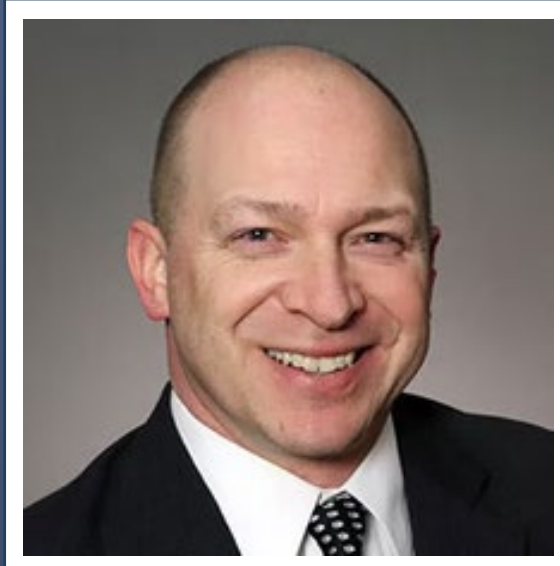
**Register:**



bit.ly/think-a-thons

- **Foster a research paradigm shift to use Big Data**
- **Promote use of Dark Data**

# Interest poll

**I am interested in (check all that apply):**

☐  Learning about Health Disparities and Health Outcomes research to apply my data science skills

☐  Conducting my own research using AI/cloud computing and publishing papers

☐  Connecting with new collaborators to conduct research using AI/cloud computing and publish papers

☐  Learning to use AI tools and cloud computing to gain new skills for research using Big Data

☐  Learning cloud computing resources to implement my own cloud

☐  Developing bias mitigation and ethical AI strategies

☐  Other

# ScHARe Guest expert

**Kenneth J. Wilkins**, PhD

NIH/NIDDK

# About Ken

Ken is a former mathematics and computer science high school teacher who found his way into biostatistics.

He worked for two decades across sectors in biomedical research, and he is now working with both NIH-employed intramural and NIH-funded extramural researchers in his NIH/NIDDK and trans-NIH roles.

His research interests encompass evolving data methods to better suit researchers' posed questions given limitations in data and data-interoperability standards.