# Preparing for AI-driven Research on ScHARe - Part 2

## A Comprehensive Review and Brainstorming Session

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

March 20, 2024

# ScHARe

**S**cience
**c**ollaborative for
**H**ealth disparities and
**A**rtificial intelligence bias
**Re**duction

**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

**NIDCR**

**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:**

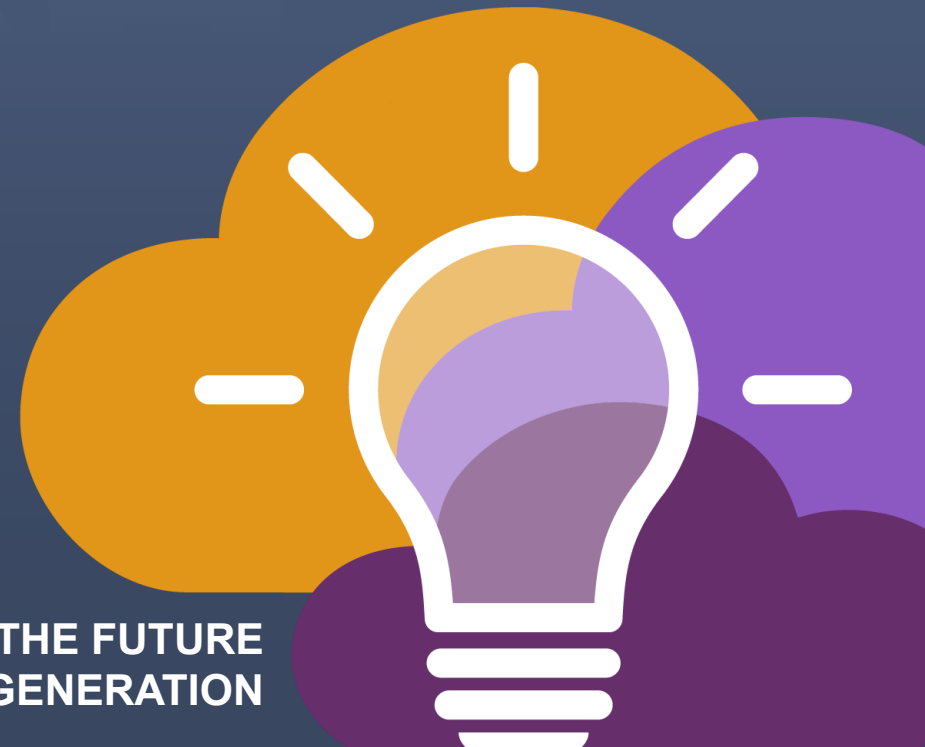|                              | None | Some | Proficient | Expert |
|------------------------------|------|------|------------|--------|
| Python                       | ☐    | ☐    | ☐          | ☐      |
| R                            | ☐    | ☐    | ☐          | ☐      |
| Cloud computing       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

# Outline

**5'**     **Introduction**
- **Experience poll**

**15'**     **ScHARe overview**
- **Interest poll**

**35'**     **Computational strategies**
- **Polls**

**20'**     **Python data science libraries**

**15'**     **Testing and monitoring in algorithm development**

**15'**     **Open science and reproducible research**

**45'**     **Research Think-a-Thons brainstorming**
- **Final poll**

**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 four 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

- Provide a **data science cloud computing resource** for community colleges and low resource minority serving institutions and organizations

# ScHARe



nimhd.nih.gov/schare

# Data Ecosystem structure
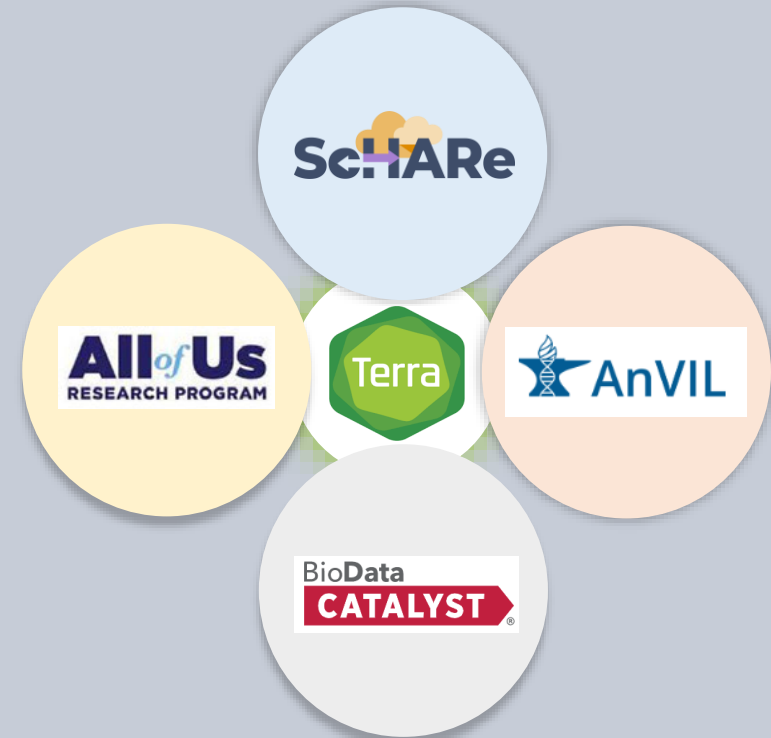## Population Science/SDoH

**240+** FEDERATED PUBLIC DATASETS
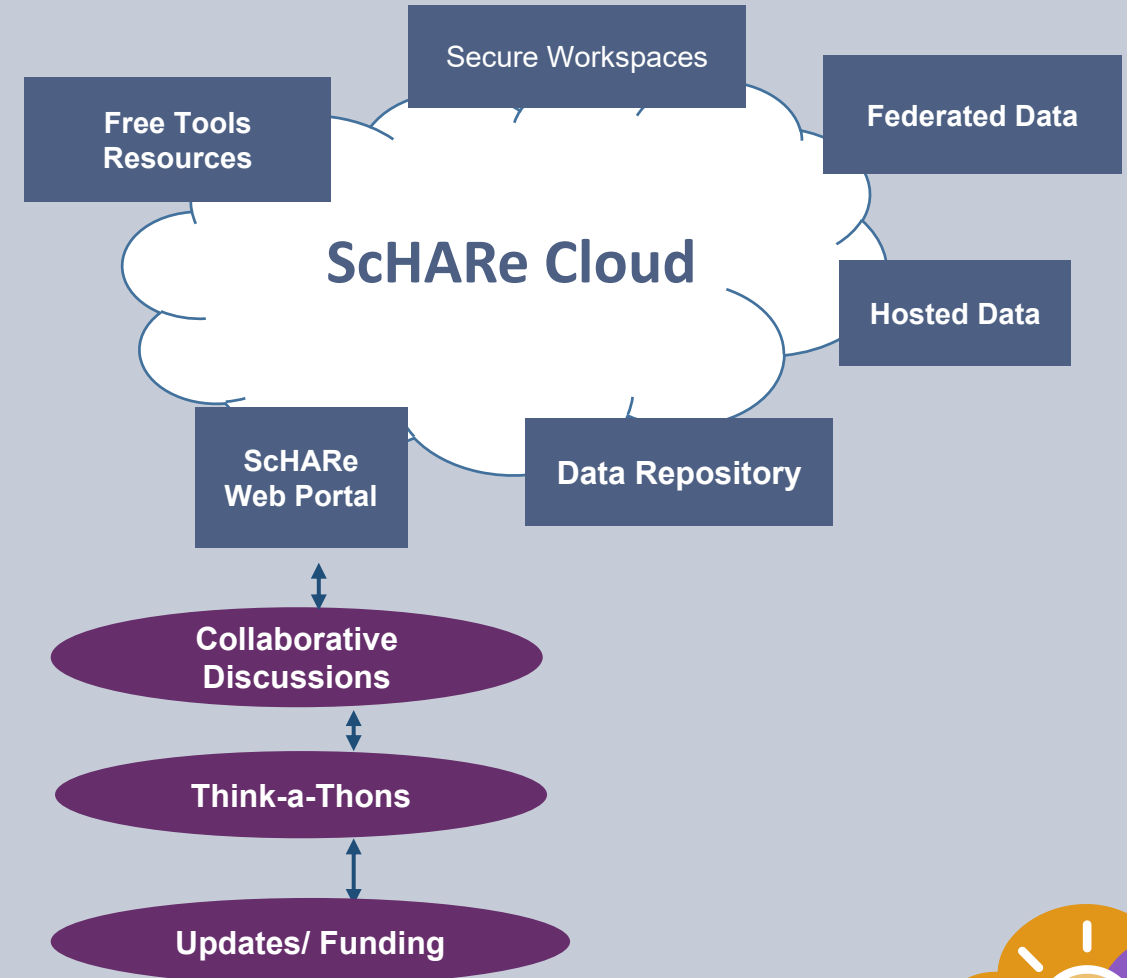
REPOSITORY

**CDE** FOCUSED

- **Population Science / SDoH / Behavioral**
- Hosted by Google & ScHARe

- **CDEs enhance data interoperability** (aggregation) by using semantic standards and concept codes

**Innovative Approach:** CDE Concept Codes Uniform Resource Identifier (**URI**)

## COMPONENTS
### Intramural and Extramural Resource

Secure Workspaces

Free Tools Resources

Federated Data

**ScHARe Cloud**

Hosted Data

ScHARe Web Portal

Data Repository

Collaborative Discussions

Think-a-Thons

Updates/ Funding

# 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 Ecosystem: Google hosted datasets

Examples of interesting datasets include:

- **American Community Survey** (U.S. Census Bureau)
- **US Census Data** (U.S. Census Bureau)
- **Area Deprivation Index** (BroadStreet)
- **GDP and Income by County** (Bureau of Economic Analysis)
- **US Inflation and Unemployment** (U.S. Bureau of Labor Statistics)
- **Quarterly Census of Employment and Wages** (U.S. Bureau of Labor Statistics)
- **Point-in-Time Homelessness Count** (U.S. Dept. of Housing and Urban Development)
- **Low Income Housing Tax Credit Program** (U.S. Dept. of Housing and Urban Development)
- **US Residential Real Estate Data** (House Canary)
- **Center for Medicare and Medicaid Services - Dual Enrollment** (U.S. Dept. of Health & Human Services)
- **Medicare** (U.S. Dept. of Health & Human Services)
- **Health Professional Shortage Areas** (U.S. Dept. of Health & Human Services)
- **CDC Births Data Summary** (Centers for Disease Control)
- **COVID-19 Data Repository by CSSE at JHU** (Johns Hopkins University)
- **COVID-19 Mobility Impact** (Geotab)
- **COVID-19 Open Data** (Google BigQuery Public Datasets Program)
- **COVID-19 Vaccination Access** (Google BigQuery Public Datasets Program)

# ScHARe Ecosystem: ScHARe hosted datasets

Organized based on the **CDC SDoH categories**, with the addition of *Health Behaviors* and *Diseases and Conditions*:

**200+** datasets

- **What are the Social Determinants of Health?**

Social determinants of health (SDoH) are the **nonmedical factors that influence health outcomes.**

They are the **conditions in which people are born, grow, work, live, and age, and the wider set of forces and systems shaping the conditions of daily life.**

Health Care and Quality

Neighborhood and Built Environment

Social and Community Context

Education Access and Quality

Economic Stability

www.cdc.gov/about/sdoh/index.html

# ScHARe Ecosystem: ScHARe hosted datasets

Examples of datasets for each category include:

## Education access and quality

Data on graduation rates, school proficiency, early childhood education programs, interventions to address developmental delays, etc.

Examples:

- **EDFacts Data Files** (U.S. Dept. of Education) - Graduation rates and participation/proficiency assessment
- **NHES - National Household Education Surveys Program** (U.S. Dept. of Education) – Educational activities

# ScHARe Ecosystem: ScHARe hosted datasets

## Health care access and quality

Data on health literacy, use of health IT, emergency room waiting times, preventive healthcare, health screenings, treatment of substance use disorders, family planning services, access to a primary care provider and high quality care, access to telehealth and electronic exchange of health information, access to health insurance, adequate oral care, adequate prenatal care, STD prevention measures, etc.

Example:

- **MEPS - Medical Expenditure Panel Survey** (AHRQ) - Cost and use of healthcare and health insurance coverage
- **Dartmouth Atlas Data -** Selected Primary Care Access and Quality Measures - Measures of primary care utilization, quality of care for diabetes, mammography, leg amputation and preventable hospitalizations

# ScHARe Ecosystem: ScHARe hosted datasets

## Neighborhood and built environment

Data on access to broadband internet, access to safe water supplies, toxic pollutants and environmental risks, air quality, blood lead levels, deaths from motor vehicle crashes, asthma and COPD cases and hospitalizations, noise exposure, smoking, mass transit use, etc.

Examples:

- **National Environmental Public Health Tracking Network** (CDC) - Environmental indicators and health, exposure, and hazard data
- **LATCH - Local Area Transportation Characteristics for Households** (U.S. Dept. of Transportation) – Local transportation characteristics for households

# ScHARe Ecosystem: ScHARe hosted datasets

## Social and community context

Data on crime rates, imprisonment, resilience to stress, experiences of racism and discrimination, etc.

Example:

- **Hate crime statistics** (FBI) - Data on crimes motivated by bias against race, gender identity, religion, disability, sexual orientation, or ethnicity
- **General Social Survey** (GSS) - Data on a wide range of characteristics, attitudes, and behaviors of Americans.

# ScHARe Ecosystem: ScHARe hosted datasets

## Economic stability

Data on unemployment, poverty, housing stability, food insecurity and hunger, work related injuries, etc.

Examples:

- **Current Population Survey (CPS) Annual Social and Economic Supplement** (U.S. Bureau of Labor Statistics ) - Labor force statistics: annual work activity, income, health insurance, and health

- **Food Access Research Atlas** (U.S. Dept. of Agriculture) – Food access indicators for low-income and other census tracts

# ScHARe Ecosystem: ScHARe hosted datasets

## Health behaviors

Data on health-related practices that can directly affect health outcomes.

Examples:

- **BRFSS - Behavioral Risk Factor Surveillance System** (CDC) - State-level data on health-related risk behaviors, chronic health conditions, and use of preventive services
- **YRBSS - Youth Risk Behavior Surveillance System** (CDC) – Health behaviors that contribute to the leading causes of death, disability, and social problems among youth and adults

# ScHARe Ecosystem: ScHARe hosted datasets

## Diseases and conditions

Data on incidence and prevalence of specific diseases and health conditions.

Examples:

- **U.S. CDI - Chronic Disease Indicators** (CDC) - 124 chronic disease indicators important to public health practice
- **UNOS - United Network of Organ Sharing** (Health Resources and Services Administration) – Organ

  transplantation: cadaveric and living donor characteristics, survival rates, waiting lists and organ disposition

# Terra Interface: Datasets and Access to Data

## Analyses

Tab in ScHARe workspace, the notebook **00_List of Datasets Available on ScHARe** lists all of the datasets available in the ScHARe Datasets collection

## Data Tab in

ScHARe workspace, **data tables help access ScHARe data and keep track of your project data:**

- ScHARe workspace, click on the Data tab
- Under Tables, see a list of dataset categories
- Click on a category, to see a list of relevant datasets
- Scroll to the right to learn more about each dataset

# Terra Interface: Secure workspace



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

- **Assign roles**: review or admin

- **Host own data and code**

# Terra Interface: Notebooks for Analytics & Tutorials

# Workflows Modular codes

A notebook integrates code and its output into a single document where you can run code, display the output, and also add explanations, formulas, and charts



**Easy to Use--Cut and Paste Analytics**



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

# ScHARe

# AGGRAGATING DATA SETS TOOL: Variables & CDEs

## CDE Mapping Project & Federated Data

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

**+** AMERICAN COMMUNITY SURVEY

**+** 

**+** **BRFSS**

**+** 

**+** **Medical Expenditure Survey (MEPS)**

***\*\*Assign DOIs to aggregated data sets***

## CDE Mapping across Program Projects

NIMHD Grantee Data Proj 1

ICO Grantee Data Proj 4

**ScHARe CDEs**

Contract Data Proj 2

Intramural Data Proj 3

## BE A PART OF THE FUTURE OF KNOWLEDGE GENERATION