# ScHARe

## Algorithm Testing and Monitoring

# What are Algorithms?

- An **algorithm** is a finite set of well-defined instructions designed to perform a specific task or solve a particular problem, often expressed in a logical sequence in a step-by-step process that can be executed by a computer

- Algorithms enable efficient and accurate **decision-making** and problem-solving across various domains, including healthcare policy decisions

# Here's how it works in healthcare

- **Lots of data:** Hospitals collect tons of data about patients (diagnoses, treatments, and meds)

- **Data analysis:** Algorithms sift through this massive amount of data and identify patterns

- **Informing decisions:** Healthcare policymakers might use these patterns to decide how to allocate resources, like which treatments are most effective or where to offer more services

Sounds helpful, right? But there's a catch:

## Algorithmic Bias

# Algorithmic Bias

- Algorithmic bias refers to systematic and unfair outcomes arising from algorithms used for decision-making

- Algorithms trained on biased data or with flawed design can perpetuate or amplify existing societal biases in healthcare

**Where can bias creep in** during the development and implementation of algorithms?

1. **Data Acquisition and Selection:**

- **Sampling Bias: i**f the data used to train the algorithm doesn't represent the entire target population

- **Historical Bias: i**f historical healthcare data reflects past discrimination

2. **Feature Engineering and Model Design:**

- **Choosing the wrong features** can lead the algorithm to make unfair decisions based on irrelevant factors.

- **Model design:** inherent limitations

3. **Model Training and Evaluation:**

- **Training data quality:** inaccurate or incomplete datas

- **Evaluation metrics: f**ocusing solely on overall accuracy might mask disparate impacts on different populations. We need fairness metrics that assess how the algorithm performs across different subgroups (e.g., race, ethnicity, socioeconomic status)

4. **Implementation and Monitoring:**

- **Limited transparency:** If the decision-making process of the algorithm is a "black box," it's hard to identify and mitigate biases

- **Unintended consequences:** Even well-intentioned algorithms can lead to unintended consequences if not continuously monitored for potential biases emerging in real-world use

# Importance of Algorithm Testing

1. Crucial during the design phase of a research project

2. Ensures the reliability and validity of algorithms before implementation

3. Enhances the accuracy and effectiveness of algorithms in real-world applications

# Importance of Algorithm Monitoring

1. **Evolving Data and Real-World Use:** The data an algorithm encounters in real-world use might differ from the training data

2. **Unforeseen Consequences:** Even well-designed algorithms can have unintended consequences

3. **Shifts in Societal Biases:** Societal biases are constantly evolving. Monitoring helps ensure the algorithm doesn't become biased due to changes in the social landscape

4. **Building Trust and Transparency:** Regular monitoring demonstrates a commitment to fairness and helps build trust in the algorithms used for healthcare decisions

# Avoiding Perpetuating Bad AI

Strategies to mitigate bias in datasets:

1. **Identify potential sources of bias:** Analyze data collection methods, sampling procedures, and variable selection for potential biases

2. **Utilize bias mitigation techniques:** Apply techniques like data balancing, weighting, or fairness-aware algorithms to mitigate bias in the data

3. **Promote transparency and responsible AI practices:** Document the limitations of the data and potential biases to ensure responsible use of AI models trained on the dataset.

# Legal and Regulatory Frameworks

Legal and regulatory frameworks govern the use of algorithms include:

1. **Anti-Discrimination Laws:** Prohibiting discrimination based on protected characteristics such as race, gender, or age

2. **Privacy Regulations:** Safeguarding individuals' privacy rights and regulating the collection and use of personal data

3. **Ethical Guidelines:** Providing guidelines for ethical algorithm development and deployment, issued by professional organizations or government agencies

**Compliance is essential**

# Ethical Considerations and Responsible AI

Principles of responsible AI include:

1. **Fairness:** Ensuring algorithms produce unbiased outcomes across different demographic groups

2. **Transparency:** Making algorithms transparent and understandable to stakeholders

3. **Accountability:** Holding developers and users accountable for the impact of algorithms

4. **Privacy:** Protecting individuals' privacy rights and sensitive information

**Adhering to ethical principles is essential for building trust and mitigating potential harms associated with AI**

# Quiz 4

**To mitigate bias in algorithms used for real-world applications, it's important to:**

a) Only use the algorithm on datasets with perfectly balanced representation

b) Continuously monitor the algorithm's performance across different demographics and adjust as needed

c) Focus solely on optimizing the accuracy of the algorithm during development

d) Limit the complexity of the algorithm to ensure easy interpretability

# Quiz 5

**Societal biases can potentially be reflected in algorithm output because:**

a) Algorithms are inherently malicious and designed to discriminate

b) Algorithms are completely objective and not influenced by external factors

c) Algorithms learn from data, which can contain societal biases

d) Algorithms are programmed by biased human creators

# Introduction to Open Science

Open Science is a paradigm shift in research practices aimed at fostering **transparency, collaboration, and accessibility**

It promotes the sharing of:
- research data
- methodologies
- findings

to accelerate scientific progress and innovation

# Open Science Principles

1.  **Open Access:** Making research findings freely available online, often through open access journals or repositories

2.  **Open Data:** Sharing the raw data used in research studies

3.  **Open Methodology:** Making the research methods and protocols used in a study openly available

4.  **Open Source:** Using and sharing open-source software for data analysis and other research tasks

5. **Open Peer Review:** Making the peer review process more transparent, allowing reviewers' identities or comments to be disclosed to some extent

6. **Reproducibility:** Conducting research in a way that allows others to reproduce the findings

7. **Collaboration:** Encouraging researchers to work together and share their findings openly

8. **Public Engagement:** Communicating scientific findings to the public in a clear, understandable way

# NIH and ScHARe Embrace Open Science

Promoting Open Science is crucial for:

- advancing **knowledge** discovery

- improving research **reproducibility**

- promoting public **trust** in science

# The new NIH Data Sharing Policy



The National Institutes of Health (NIH) implemented a new Data Management and Sharing (DMS) Policy in January 2023

**Goal:** promote transparency and responsible data management in scientific research

# The new NIH Data Sharing Policy

**Who is affected:**

- Researchers applying for NIH funding (grants, contracts)
- Intramural NIH researchers (conducting research within the NIH itself)

- **Core principle:** Maximize the appropriate sharing of scientific data

# The new NIH Data Sharing Policy

## What data needs to be shared?

- Scientific data generated from the funded/conducted research, with exceptions for data with privacy risks, commercialization potential, or security concerns

## How is data shared?

- Researchers must submit a *Data Management and Sharing Plan* outlining how they will handle data, ensure its quality and security, and deposit it in a suitable public repository

The **ScHARe** repository is designed to meet the data sharing requirements of the NIH data sharing policy

# The **ScHARe** Repository for Data Management:

- serves as a **centralized platform** for storing, managing, and sharing research data related to health disparities

- adheres to **FAIR principles** (Findable, Accessible, Interoperable, Reusable) to ensure data discoverability and usability

- supports **Open Science** initiatives and promotes collaborative research

The ScHARe repository focuses on **Common Data Elements** (CDEs), standardizing data and metadata to facilitate interoperability and data reuse

# The ScHARe CDEs

- **Common Data Elements** (CDEs) are standardized, precisely defined data points used consistently across different research studies

- They act as building blocks for collecting and sharing data in a comparable and interoperable manner

NIH CDE Repository

# Benefits of CDEs

- **Standardization:** Ensures consistency in data collection, formatting, and reporting across studies

- **Interoperability:** Facilitates data integration and comparison between different datasets and projects

- **Efficiency:** Streamlines data management processes, reducing redundancy and errors in data handling

- **Collaboration:** Promotes collaboration and data sharing among researchers

- **Quality:** Enhances data quality and reliability by adopting standardized data collection and reporting practices

# ScHARe Core CDEs

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

**NIH Endorsed**

- 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 (and Associated Medications/Treatments)

- **NIMHD Framework***
- **Health Disparity Outcomes***

\* Project Level CDEs

**For FUNDED PROJECT DATA** – CDEs Centralized for Interoperability and Data Sharing

# Quiz 6

**The core principle of the NIH Data Management and Sharing (DMS) Policy emphasizes:**

a) Restricting access to all scientific data generated by NIH-funded research

b) Maximizing the appropriate sharing of scientific data while considering ethical and legal limitations

c) Encouraging the publication of research findings in open-access journals only

d) Prioritizing data privacy over all other considerations

# Quiz 7

**Which of the following is a primary benefit of using common data elements (e.g., standardized variable names, units of measure)?**

a) Improves data security and privacy

b) Simplifies data analysis and comparison across studies

c) Reduces the overall size of data storage requirements

d) Enhances the visual appeal of data presentations

# Let's brainstorm health disparities research ideas

**Let's consider:**

- **innovative approaches and methodologies, such as AI**

- **datasets publicly available on ScHARe**

# Health disparities

A health disparity is a health difference that adversely affects disadvantaged **populations** in comparison to a reference population, based on one or more **health outcomes**

## Health Disparity Outcomes

The health outcomes are categorized as:

- Higher incidence and/or prevalence of disease, including earlier onset or more aggressive progression of disease.
- Premature or excessive mortality from specific health conditions.
- Greater global burden of disease, such as Disability Adjusted Life Years (DALY), as measured by population health metrics.
- Poorer health behaviors and clinical outcomes related to the aforementioned.
- Worse outcomes on validated self-reported measures that reflect daily functioning or symptoms from specific conditions.

## Populations with Health Disparities

Populations that experience health disparities include:

- Racial and ethnic minority groups
- People with lower socioeconomic status (SES)
- Underserved rural communities
- Sexual and gender minority (SGM) groups
- People with disabilities

# Inequities can lead to health disparities

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**

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

**Health Care and Quality**

**Neighborhood and Built Environment**

**Social and Community Context**

**Education Access and Quality**

**Economic Stability**

If certain communities have less access to good education, jobs, fresh food or healthcare, they might face **more challenges in staying healthy** or may not have the same **opportunities to make healthy choices**

How do these **nonmedical factors interact with each other and biology** to influence health?

**Artificial Intelligence** may have the answer

# Identifying research gaps

Areas with limited research in the current health disparity research landscape:

1. **Social Determinants of Health (SDoH) Interactions**

   A gap exists in understanding the complex interactions between SDoH factors and how they contribute to health disparities across different populations

2. **Precision Disparities**

   The rise of personalized medicine using genomics raises concerns about potential disparities in access and benefits. How do genetic and social factors intertwine to create "precision health disparities"?

3. **Intersectionality and Health**
   Traditional research focuses on single demographic factors (e.g., race, gender). How do multiple social identities (e.g., Black woman, LGBTQ+, immigrant) intersect and influence health outcomes?

4. **Role of Implicit Bias in Healthcare Systems**
   How does implicit bias in healthcare delivery affect treatment recommendations and patient experiences? What interventions can mitigate its effects?

5. **Digital Divide and Disparities**

   Lack of access to technology can exacerbate health disparities. How to leverage technology to improve health outcomes for underserved populations while ensuring equitable access?

6. **Environmental Exposures and Disparities**
   Communities of color and low-income populations are often disproportionately exposed to environmental hazards. What are the long-term health effects of these exposures?

7. **Longitudinal Studies on Disparities**
   Many studies are cross-sectional. Longitudinal studies that track individuals over time are crucial for understanding the progression of disparities

What data is available?

**The ScHARe Data Ecosystem**

# SDoH-related Datasets Available on ScHARe: A Valuable Resource

ScHARe provides a valuable platform for researchers seeking **SDoH-related data**

**Explore the available datasets** to identify potential resources that align with your research interests in social determinants of health and their impact on various health outcomes

# ScHARe Ecosystem

**The ScHARe Data Ecosystem is comprised of:**

1. **Google Hosted Public Datasets:** publicly accessible, federated, de-identified datasets hosted by Google through the Google Cloud Public Dataset Program
   **Example**: *American Community Survey (ACS)*

2. **ScHARe Hosted Public Datasets:** publicly accessible, de-identified datasets hosted by ScHARe
   **Example**: *Behavioral Risk Factor Surveillance System (BRFSS)*

3. **ScHARe Hosted Project Datasets:** 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*

# 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*:

**240+** 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

# How to check what data is available on ScHARe

## Analyses tab

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

ScHARe datasets

# The ScHARe Data Ecosystem

Last updated: November 27, 2023

This document is intended to provide a comprehensive list of the datasets available in the ScHARe Data Ecosystem for analysis in the ScHARe Terra instance. Using the ScHARe Data Ecosystem, researchers are able to search, link, share, and contribute to a collection of datasets relevant to social science, health outcomes, minority health and health disparities research.

The collection is comprised of:

- **Google-hosted Public Datasets** - Publicly accessible, federated, de-identified datasets hosted by Google through the Google Cloud Public Dataset Program. Examples: US Census Data; American Community Survey (ACS)
- **ScHARe-hosted Public Datasets** - Publicly accessible, de-identified datasets hosted by ScHARe. Examples: Social Vulnerability Index (SVI), Behavioral Risk Factor Surveillance System (BRFSS)
- **ScHARe-hosted Project Datasets** - Publicly accessible and controlled-access, funded program/project datasets shared by NIH grantees and intramural investigators to comply with the program/project datasets shared by NIH grantees and intramural investigators to comply with the Jackson Heart Study (JHS)

Scan me

bit.ly/ScHARe-datasets

# How to access available data on ScHARe

## Data tab

In the **Data** tab in the ScHARe workspace, **data tables help access ScHARe data and keep track of your project data:**

- In the ScHARe workspace, click on the Data tab

- Under Tables, you will see a list of dataset categories

- If you click on a category, you will see a list of relevant datasets

- Scroll to the right to learn more about each dataset

# Potential projects leveraging AI and Big Data

These examples showcase the intersection of different **datasets** and the application of diverse **AI tools** to gain insights into the social determinants of health and their impact on health outcomes and disparities in minority populations

# #1: Geospatial Analysis of Environmental Factors

- **Objective:** Explore the impact of environmental conditions on health outcomes, especially in minority communities

- **Methodology:**
  - Combine environmental datasets (air quality, pollution levels) with health records using geospatial analytics
  - AI models can reveal spatial patterns, helping identify areas with higher health risks in minority populations

This information can inform policies addressing environmental justice and public health

# #1: Geospatial Analysis of Environmental Factors

- **Datasets:**
  - Environmental Protection Agency (EPA) Air Quality Data: Provides information on air pollutants and air quality indices
  - Health and Nutrition Examination Survey (NHANES): Includes health data, including respiratory health indicators

- **AI Tools:**
  - Geospatial Analytics Tools: Geographic Information System (GIS) platforms like ArcGIS or QGIS to map environmental data and health outcomes
  - Machine Learning for Spatial Analysis: Algorithms for spatial regression or clustering to identify areas with higher health risks

# #2: Education and Health Disparities Analysis

- **Objective:** Examine the link between educational disparities and health outcomes in minority communities

- **Methodology:**
  - Merge educational attainment data with health records, applying AI techniques to discern patterns
  - Explore how educational opportunities influence health behaviors, preventive care, and overall well-being

This interdisciplinary research can inform education and public health policies aimed at reducing health disparities

# #2: Education and Health Disparities Analysis

- **Datasets:**
  - National Center for Education Statistics (NCES) Educational Attainment Data: Contains data on educational attainment by demographics
  - Behavioral Risk Factor Surveillance System (BRFSS): Includes self-reported health data and behaviors

- **AI Tools:**
  - Predictive Modeling: Utilize algorithms like logistic regression or neural networks to predict health outcomes based on educational disparities
  - Causal Inference Techniques: Apply methods such as propensity score matching to isolate the impact of education on health

# #3: Causal links between chronic stress associated with social adversity and health disparities

**Health is adversely affected by social disadvantage**:

1. **Neighborhoods influence health through their physical and geographic characteristics**:
   - air and water quality
   - lead paint exposure
   - proximity to health promoting features (e.g.: hospitals, healthy food stores)
   - proximity to health suppressing features (e.g.: toxic factories, fast food)
   - access to green space, etc.

2. **Chronic stress of social disadvantage, socioeconomic inequality, and racial discrimination can influence health** through a variety of biological pathways:
   1. neuroendocrine
   2. developmental
   3. immunologic
   4. vascular

**Objective:** Examine the role of epigenetic modifications as a causal link between chronic stress associated with social adversity and health disparities, and impact of mitigating factors

# Research Projects Brainstorming

**What research projects do you believe it would be worthwhile to pursue?**

# ScHARe resources

Support made available to users:

**ScHARe-specific**
- ScHARe documentation
- Email support

**Platform-specific**
- Terra-specific support
- Terra-specific documentation

# ScHARe resources

Training opportunities made available to users:

- Monthly **Think-a-Thons**

- **Instructional materials** and slides made available online on NIMHD website

- **YouTube videos**

- **Links to relevant online resources** and training on NIMHD website

- **Pilot credits** for testing ScHARe for research needs

- **Instructional Notebooks** in ScHARe Workspace with instructions for:

    - Exploring the data ecosystem

    - Setting your workspace up for use

    - Accessing and interacting with the categories of data accessible through ScHARe

# ScHARe resources: cheatsheets



Credits: datacamp.com

# Terra resources

If you are new to Terra, we recommend exploring the following resources:

- Overview Articles: Review high-level docs that outline what you can do in Terra, how to set up an account and account billing, and how to access, manage, and analyze data in the cloud
- Video Guides: Watch live demos of the Terra platform's useful features
- Terra Courses: Learn about Terra with free modules on the Leanpub online learning platform
- Data Tables QuickStart Tutorial: Learn what data tables are and how to create, modify, and use them in analyses
- Notebooks QuickStart Tutorial: Learn how to access and visualize data using a notebook
- Machine Learning Advanced Tutorial: Learn how Terra can support machine learning-based analysis

# ScHARe

Thank you

# Think-a-Thon poll

1. **Rate how useful this session was:**

☐ Very useful

☐ Useful

☐ Somewhat useful

☐ Not at all useful

# Think-a-Thon poll

2.    Rate the pace of the instruction for yourself:

☐ Too fast

☐ Adequate for me

☐ Too slow

# Think-a-Thon poll

3. How likely will you participate in the next Think-a-Thon?

☐ Very interested, will definitely attend

☐ Interested, likely will attend

☐ Interested, but not available

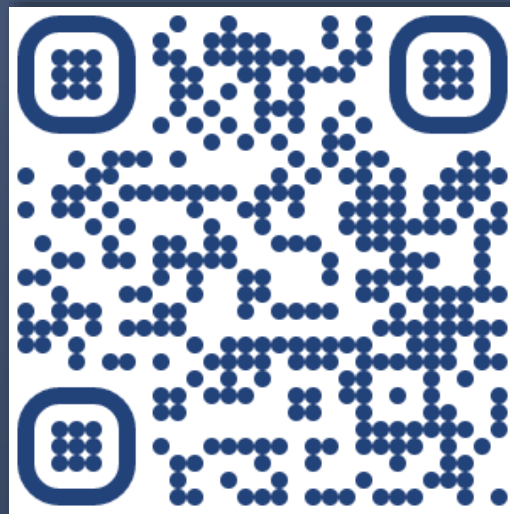☐ Not interested in attending any others

# ScHARe

**Next Think-a-Thons:**

**Register for ScHARe:**

✉ schare@mail.nih.gov

bit.ly/think-a-thons

bit.ly/join-schare