universidade de aveiro **sbidm** serviços de biblioteca, informação documental e museologia
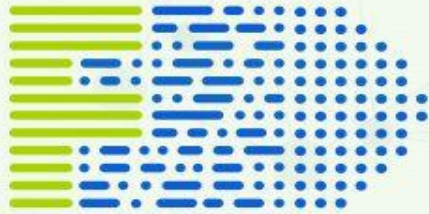
**AMERICAN CORNERS PORTUGAL**

**22 DE NOVEMBRO 2019**
**UNIVERSIDADE DE AVEIRO**

**5.º FÓRUM**
# GESTÃO DE DADOS DE INVESTIGAÇÃO

RCAAP | Repositórios Científicos de Acesso Aberto de Portugal

Organização

REPÚBLICA PORTUGUESA
CIÊNCIA, TECNOLOGIA E ENSINO SUPERIOR

**FCT** Fundação para a Ciência e a Tecnologia

Apoio

universidade de aveiro

**AMERICAN CORNERS PORTUGAL**

ifdo
INTERNATIONAL FEDERATION OF DATA ORGANIZATION

GDCC

**THE ODUM INSTITUTE**
FOR RESEARCH IN SOCIAL SCIENCE

# Research Data Curation, Management, Sharing and Archiving

Jonathan Crabtree

Director of Cyberinfrastructure

Odum Institute for Research in Social Science

# THE H. W. ODUM INSTITUTE FOR RESEARCH IN SOCIAL SCIENCE

- Founded in 1924 by Howard W. Odum

- Oldest university-based interdisciplinary social science research institute in the U.S.

- *To help grow and lead a world-class social science research infrastructure at the University of North Carolina at Chapel Hill to ensure that researchers can conduct scientifically rigorous research that contributes to better lives of the citizens of North Carolina and the World…*

**THE ODUM INSTITUTE**
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
*of* NORTH CAROLINA
*at* CHAPEL HILL

# Research Data Management: Context is Critical

What are the components of your research data environment?

What would you need to reproduce your own research?

**THE ODUM INSTITUTE**
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# GETTING TO KNOW (AND LOVE) RESEARCH DATA

What are Data?

Research Data Management

Data Documentation

Data Security

Archiving Data with Dataverse

Odum Institute Services and Ongoing Projects

www.digitalbevaring.dk

**THE ODUM INSTITUTE**
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
*of* NORTH CAROLINA
*at* CHAPEL HILL

DATA

The real world

Idea of the real world

What can be measured

What can be learned from analysis

Little, J., & Zoss, A. (2014, September). *Basic data cleaning and analysis for data tables.* Webinar, Duke University Perkins Library. Retrieved from http://library.capture.duke.edu/Panopto/Pages/Viewer.aspx?id=9e7b8529-3566-4469-98f3-4e520f32b849

THE ODUM INSTITUTE
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# WHAT ARE DATA?

**The world as DATA**

- People

- Objects

- Places/Spaces

- Time

- Relationships

- Ideas/Concepts

**THE ODUM INSTITUTE**
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# WHAT ARE DATA?

- Telephone interview of 2,002 adults 18 or older
- Randomly selected youngest adult in household



Adapted photo by James Cridland available under CC BY 2.0 at
https://flic.kr/p/Wd54U

THE ODUM INSTITUTE

FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# Research Data Formats are Evolving

Increasing in size

Becoming more dynamic

More diverse

Harder to de-identify

More difficult to integrate with other data sources

THE ODUM INSTITUTE
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# Curation Processes

How do curation processes differ across file formats and disciplinary contexts?

THE ODUM INSTITUTE
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# WHY DATA MANAGEMENT?

Data management refers to activities that support long-term preservation, access, and use of data.

- Planning for data management

- Describing data

- Formatting data

- Storing and backing up data

- Anonymizing data

- Controlling access to data

THE ODUM INSTITUTE
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# WHY DATA MANAGEMENT?

- Data management makes it possible for other researchers to discover, interpret, and re-use data.

- Data management helps sustain the value of data by enabling others to verify and build upon published results.

- Data management facilitates long-term preservation of and access to data.

**THE ODUM INSTITUTE**
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
*of* NORTH CAROLINA
*at* CHAPEL HILL

# DATA MANAGEMENT PLANNING

Data management is most successful when data management practices are implemented throughout the research lifecycle.



Source: University of Virginia Library. (2014). Steps in the research lifecycle. Retrieved September 21, 2014, from http://dmconsult.library.virginia.edu/lifecycle/

**THE ODUM INSTITUTE**
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# DATA DOCUMENTATION

…**sufficient information** exists with which to understand, evaluate, and build upon a prior work if a third party could replicate the results **without any additional information from the author**.

King, G. (1995). Replication, replication. *PS: Political Science & Politics*, *28*(3), 444–452. https://doi.org/10.2307/420301

Image source: http://harvardmagazine.com/2009 /09/two-honored-with-university-professorships

# DATA DOCUMENTATION



**CODEBOOK**

**README**

**ANALYSIS CODE**

THE ODUM INSTITUTE
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# DATA DOCUMENTATION

**CODEBOOK**

- Variable names + labels
- Value codes + labels
- Range of values
- Data type

THE ODUM INSTITUTE
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# DATA DOCUMENTATION

**README**

- Data collection methods
- Coding information
- Variable construction
- Dataset modifications
- Original data source

THE ODUM INSTITUTE
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# DATA DOCUMENTATION

- Software version
- Commands
- Comment statements

**ANALYSIS CODE**

THE ODUM INSTITUTE
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# DATA SECURITY

- Names
- Geographic subdivisions smaller than state
- Zip codes
- All elements of dates except year directly related to an individual
- Telephone numbers
- Fax numbers
- Email addresses
- Social Security numbers
- Medical record numbers
- Health plan beneficiary identifiers
-  Account numbers
- Certificate/license numbers
- Vehicle identifiers and serial numbers
- Device identifiers and serial numbers
- Web universal resource locators (URL)
- Internet protocol (IP) address numbers
- Biometric identifiers
- Full face photographic images
- Any other number, characteristic, or code that could be used by the researcher to identify the individual

**THE ODUM INSTITUTE**
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# DATA SECURITY

http://aboutmyinfo.org/

Date of birth
Gender
Zip code

$+$

**87%**

uniquely
identifiable

Sweeney, L. (2000). *Simple demographics often identify people uniquely* (Data Privacy Working Paper No. 3). Pittsburgh, PA: Carnegie Mellon University. Retrieved from http://dataprivacylab.org/projects/identifiability/paper1.pdf

THE ODUM INSTITUTE
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# Data security: threats and vulnerabilities

Sources of threat
- Natural
- Unintentional Human
- Intentional

Areas of vulnerability
- **Logical**: Data at rest in system, data in motion across networks, data being processed in applications
- **Physical**: Computer systems, network, and backups, disposal, media
- **Social**: social engineering, mistakes, insider threats

Altman, Micah. (2013). *Managing Confidential Data* [PowerPoint slides]. Retrieved from http://www.slideshare.net/drmaltman/altman-confidentialdata-v22mit?ref=http://informatics.mit.edu/classes/managing-confidential-data

THE ODUM INSTITUTE
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# Confidentiality & security across the research lifecycle

**Identify** potentially sensitive information

**Review** applicable laws

**Design** risk mitigation across data lifecycle

**Reduce** sensitivity of collected data

**Plan** for publication, dissemination, and reuse

**Describe** reuse plan in consent form

**THE ODUM INSTITUTE**
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
*of* NORTH CAROLINA
*at* CHAPEL HILL

# Confidentiality & security across the research lifecycle

**Separate** sensitive information in *collection*

**Encrypt** sensitive information in *transit*

**Follow** data security best practices



PLAN

CREATE

USE

ARCHIVE

# Confidentiality & security across the research lifecycle

**Protect** sensitive information in *systems*

**Desensitize** information in *processing*

**Monitor** threats and vulnerabilities

**Implement** strategies for limiting disclosure risks

**Review** data for sensitive information prior to ingest into repository

PLAN

CREATE

USE

ARCHIVE

THE ODUM INSTITUTE
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# Confidentiality & security across the research lifecycle

**Deposit** data in a trusted repository

**Dispose** of confidential data following best practices



PLAN

CREATE

USE

ARCHIVE

**THE ODUM INSTITUTE**
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
*of* NORTH CAROLINA
*at* CHAPEL HILL

# Sharing safely: New approaches

Synthetic data

Differential Privacy

Database cryptography

Rules based data sharing tools

**THE ODUM INSTITUTE**
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# THE **FAIR** DATA PRINCIPLES

FORCE11. (2016). *Guiding principles for findable, accessible, interoperable and reusable data* (Publishing Version No. B1.0). Retrieved from https://www.force11.org/fairprinciples



**F**indable  **A**ccessible  **I**nteroperable  **R**eusable

# ARCHIVING DATA

A trusted digital repository is one whose mission is to provide reliable long-term access to managed digital resources to its designated community, now and in the future.

RLG/OCLC Working Group on Digital Archive Attributes. (2002). *Trusted digital repositories: Attributes and responsibilities* (An RLG-OCLC Report). Mountain View, CA: Research Libraries Group. Retrieved from http://www.oclc.org/research/activities/past/rlg/trustedrep/repositories.pdf

www.digitalbevaring.dk

**THE ODUM INSTITUTE**
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
*of* NORTH CAROLINA
*at* CHAPEL HILL

# TRUST Principles

**FAIR** defines the properties of <u>data and metadata</u>

**TRUST** describes the characteristics of <u>data repositories</u> that are responsible for managing and disseminating the data over a long period of time

**FAIR** data in repositories we **TRUST**



**T - Transparency** is achieved by providing publicly accessible evidence of the services that a repository can and can not offer.

**R - Responsibility** is a commitment to provide high (technical) quality data services.

**U - User community** is the focus on the uses and potential uses of the data and services offered.

**S - Sustainability** is the capability to support long-term data preservation and use.

**T - Technology** is the infrastructure and capabilities to support the repository operations.

**THE ODUM INSTITUTE**
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
*of* NORTH CAROLINA
*at* CHAPEL HILL

# FINDING A REPOSITORY

1. Is the repository reputable?

2. Will it take the data you want to deposit?

3. Will it be safe in legal terms?

4. Will the repository sustain the data value?

5. Will it support analysis and track data usage?

Whyte, A. (2015). Where to keep research data: DCC Checklist for evaluating data repositories (v.1.1). Edinburgh: Digital Curation Centre. http://www.dcc.ac.uk/resources/how-guides-checklists/where-keep-research-data

THE ODUM INSTITUTE
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# ARCHIVING DATA

Include documentation and metadata
Provide information to enable discovery and appropriate interpretation and reuse of the data

README FILE

CODEBOOK

METADATA

THE ODUM INSTITUTE
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# ARCHIVING DATA

Include documentation and metadata
Provide information to enable discovery and appropriate interpretation and reuse of the data

ANALYSIS CODE

THE ODUM INSTITUTE
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# TIPS FOR ARCHIVING DATA

Archive data in open file formats
Use formats that support preservation, accessibility, and reuse of data

**DOCUMENTATION**

**DATA FILES**

**PDFA** Adobe

**.csv**
**.tab**
**.txt**

# TIPS FOR ARCHIVING DATA

Outline terms of use & apply a standard license
Enable informed reuse by clearly outlining how data can be accessed, used, and disseminated

DATA USE
AGREEMENTS

EMBARGOES

CREATIVE COMMONS
LICENSES

THE ODUM INSTITUTE
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# TIPS FOR ARCHIVING DATA

**Resolve data ownership & sharing issues early**
Discuss data sharing & archiving with collaborators, participants, and other stakeholders early in a project

THE ODUM INSTITUTE
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
*of* NORTH CAROLINA
*at* CHAPEL HILL

# Data citation

Provides stable access to data

Encourages acknowledgement and credit for data producers

Provides incentives for sharing data

Helios Herrera; Massimo Morelli; Salvatore Nunnari, 2015, "Replication data for: Turnout Across Democracies", http://dx.doi.org/10.7910/DVN/9TPNOT, Harvard Dataverse, V1 [UNF:6:NhH3rblEwGkIIbw9mygwCQ==]

# Data citation

Joint Declaration of Data Citation Principles

Synthesis of a number of groups and sponsored by Force 11

Facilitate the creation of citation practices that are both human understandable and machine-actionable

| 8 Data Citation Principles |
| --- |
| 1. Importance |
| 2. Credit & Attribution |
| 3. Evidence |
| 4. Unique Identification |
| 5. Access |
| 6. Persistence |
| 7. Specificity & Verifiability |
| 8. Interoperability & Flexibility |

Data Citation Synthesis Group: **Joint Declaration of Data Citation Principles**. Martone M. (ed.) San Diego CA: FORCE11; 2014 [https://www.force11.org/datacitation].

**THE ODUM INSTITUTE**
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
*of* NORTH CAROLINA
*at* CHAPEL HILL

# Data citation

Principle 2: Credit & Attribution

**Helios Herrera**; Massimo Morelli; Salvatore Nunnari, 2015, "Replication data for: Turnout Across Democracies", **http://dx.doi.org/10.7910/DVN/9TPNOT**, **Harvard Dataverse**, **V1** **[UNF:6:NhH3rblEwGkIIbw9mygwCQ==]**

Principle 7: Specificity & Verification (e.g. the specific version used)

Principle 4: Unique Identifier (DOI). Principle 5 & 6: Access, Persistence (A persistent identifier that provides access and metadata)

# Data citation

Assign persistent identifiers (DOIs) to data

Supports simple & effective methods of data citation, discovery, and access

Ensures data can be located online

**THE ODUM INSTITUTE**
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# The Data Sharing Problem

Different needs for archives, data libraries, researchers, journals, funding agencies…

We should preserve the data

No publications without data

We need persistent links

I want credit for my data

I need a Data Management Plan

Cross, M. *Why the Dataverse Network?*
Available at: thedata.org

# Odum's Solution

Dataverse: **centralized** professional archiving with **distributed** control and recognition

- Persistent identifiers
- Fixity
- Backups & recovery
- Metadata standards
- Conversion standards
- Preservation standards

**+**

- Branding & visibility
- Data discovery
- Ease of use
- Scholarly citation
- Control over updates
- Terms of access  & use

Cross, M. *Why the Dataverse Network?* Available at: thedata.org

THE ODUM INSTITUTE
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
*of* NORTH CAROLINA
*at* CHAPEL HILL

DATAVERSE REPOSITORIES - A WORLD VIEW

26 Installations    2,357 Dataverses    49,698 Datasets    2,747,888 Downloads

Stats generated: 11th October 2017 07:36 EDT

THE ODUM INSTITUTE
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# THE DATAVERSE PROJECT



https://dataverse.unc.edu/

# THE DATAVERSE PROJECT

- Open source web application for publishing, citing, analyzing, and preserving research data

- Data sharing and archiving with control and recognition for data producers

- Rich data support for certain file formats

- Supports data management standards and best practices

**THE ODUM INSTITUTE**
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
*of* NORTH CAROLINA
*at* CHAPEL HILL

# THE DATAVERSE PROJECT



## Archival Record

- Standardized DDI metadata

- Formal citation

- Persistent identification

https://dataverse.unc.edu/

**THE ODUM INSTITUTE**
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# External Analysis Tools via Dataverse API

# Moving beyond social science

Dataverse Network is cross-disciplinary.

We are expanding the study metadata and building communities of interested groups:
- dataverse-community@googlegroups.com

Cross, M. *Why the Dataverse Network?* Available at: thedata.org

# Support for metadata sharing

**THE ODUM INSTITUTE**
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
*of* NORTH CAROLINA
*at* CHAPEL HILL

# THE H. W. ODUM INSTITUTE FOR RESEARCH IN SOCIAL SCIENCE

**SERVICES ACROSS THE RESEARCH LIFECYCLE**

**DATA ARCHIVE**

PLAN

CREATE

USE

ARCHIVE

THE ODUM INSTITUTE FOR RESEARCH IN SOCIAL SCIENCE

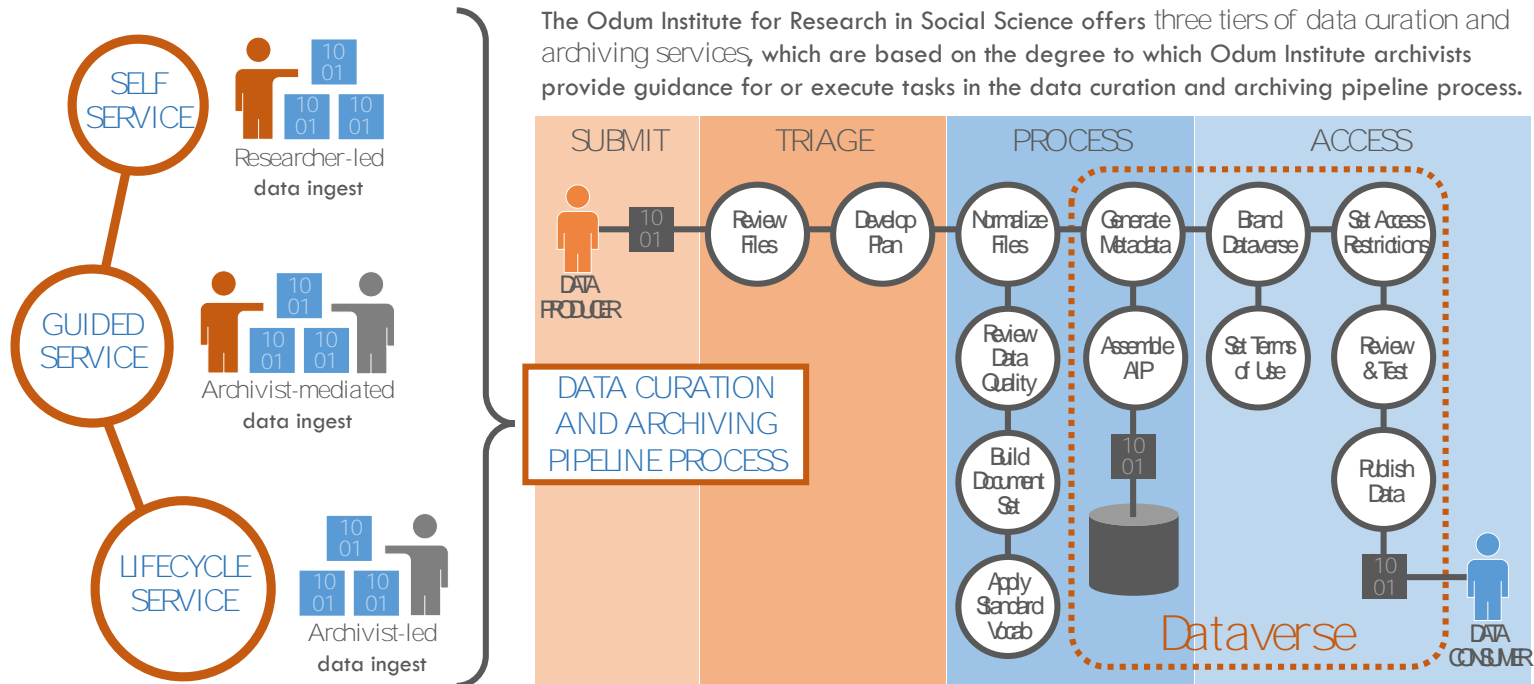THE UNIVERSITY of NORTH CAROLINA at CHAPEL HILL

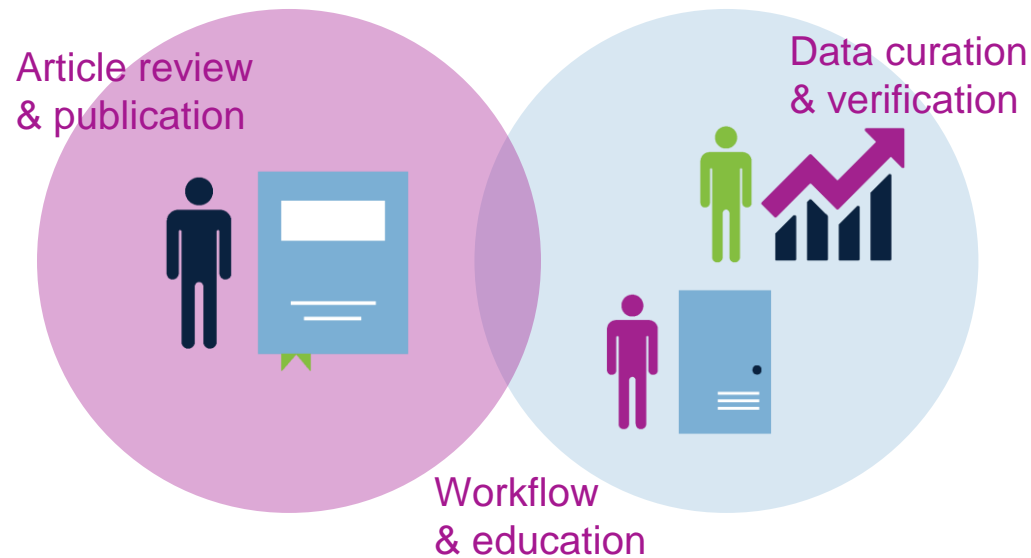# DATA CURATION AND ARCHIVING SERVICES POWERED BY **Dataverse**

The Odum Institute for Research in Social Science offers three tiers of data curation and archiving services, which are based on the degree to which Odum Institute archivists provide guidance for or execute tasks in the data curation and archiving pipeline process.

SELF SERVICE — Researcher-led **data ingest**

GUIDED SERVICE — Archivist-mediated **data ingest**

LIFECYCLE SERVICE — Archivist-led **data ingest**

DATA CURATION AND ARCHIVING PIPELINE PROCESS

**SUBMIT**
DATA PRODUCER
- Review Files
- Develop Plan

**TRIAGE**

**PROCESS**
- Normalize Files
- Review Data Quality
- Build Document Set
- Apply Standard Vocab
- Generate Metadata
- Assemble AIP

**ACCESS**
- Brand Dataverse
- Set Access Restrictions
- Set Terms of Use
- Review & Test
- Publish Data

Dataverse

DATA CONSUMER

# Data Verification Service



Article review & publication

Data curation & verification

Workflow & education

THE ODUM INSTITUTE
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# Data verification

THE ODUM INSTITUTE
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# Confirmable Reproducible Research (CoRe2) Environment

# Impact Project Overview

- **Challenges:**
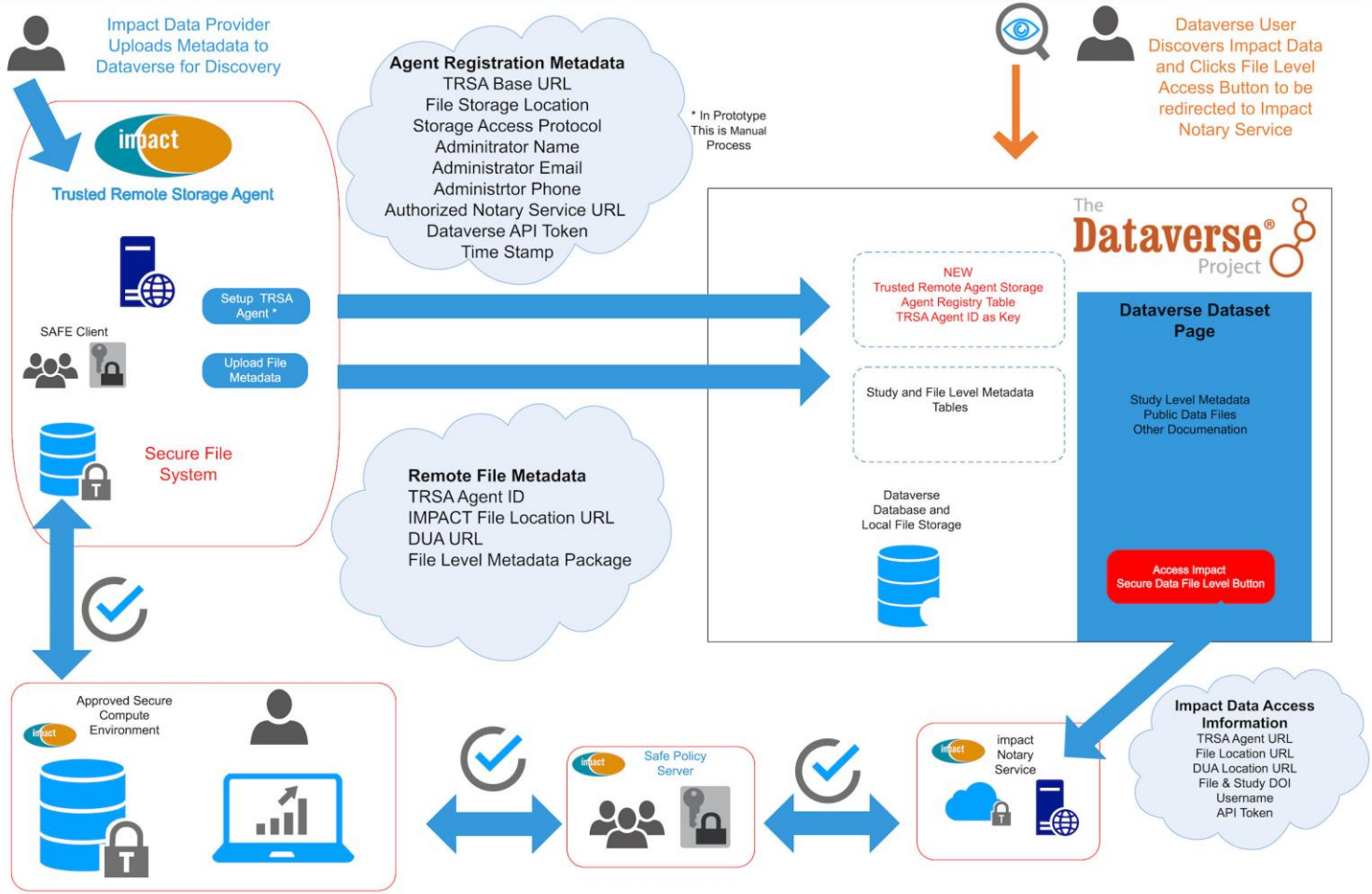  - Social science and many other data-oriented disciplines depend on data belonging to multiple stakeholders
  - Governed by a variety of use policies
  - Multi-institutional research requires cooperative analysis
  - Need to satisfy the privacy concerns of the owners while producing interesting research outcomes by analyzing data

- **Goal:** to enable cooperative processing across the stakeholder-owned datasets, while respecting the privacy policies of the individual owners, <u>and</u> to provide a model for collaboration that could be readily used by other institutions.

**THE ODUM INSTITUTE**
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

IMPACT TRSA User Workflow

THE ODUM INSTITUTE
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

# Thank You

## CONNECT WITH
## THE ODUM INSTITUTE

Jonathan Crabtree
Jonathan_Crabtree@unc.edu

**The Odum Institute**
http://www.odum.unc.edu

@Odum_Institute

The Odum Institute

THE ODUM INSTITUTE
FOR RESEARCH IN SOCIAL SCIENCE

www.odum.unc.edu

THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL