# Principles of Transparent Research
## *Implementation Challenges*

Lars Vilhuber

Cornell University
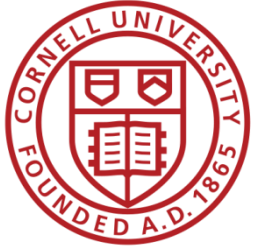
# Context

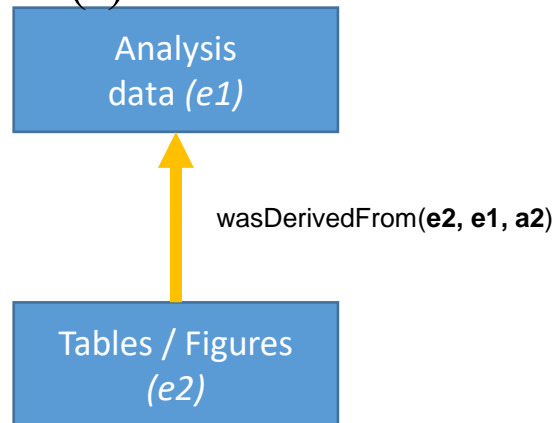wasDerivedFrom(e2, e1, a, g2, u1)

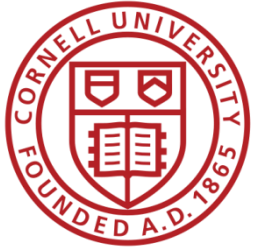# Replication packages in Social Sciences

- Started to appear in mid-1980s
- American Economic Review "data availability policy" 2005
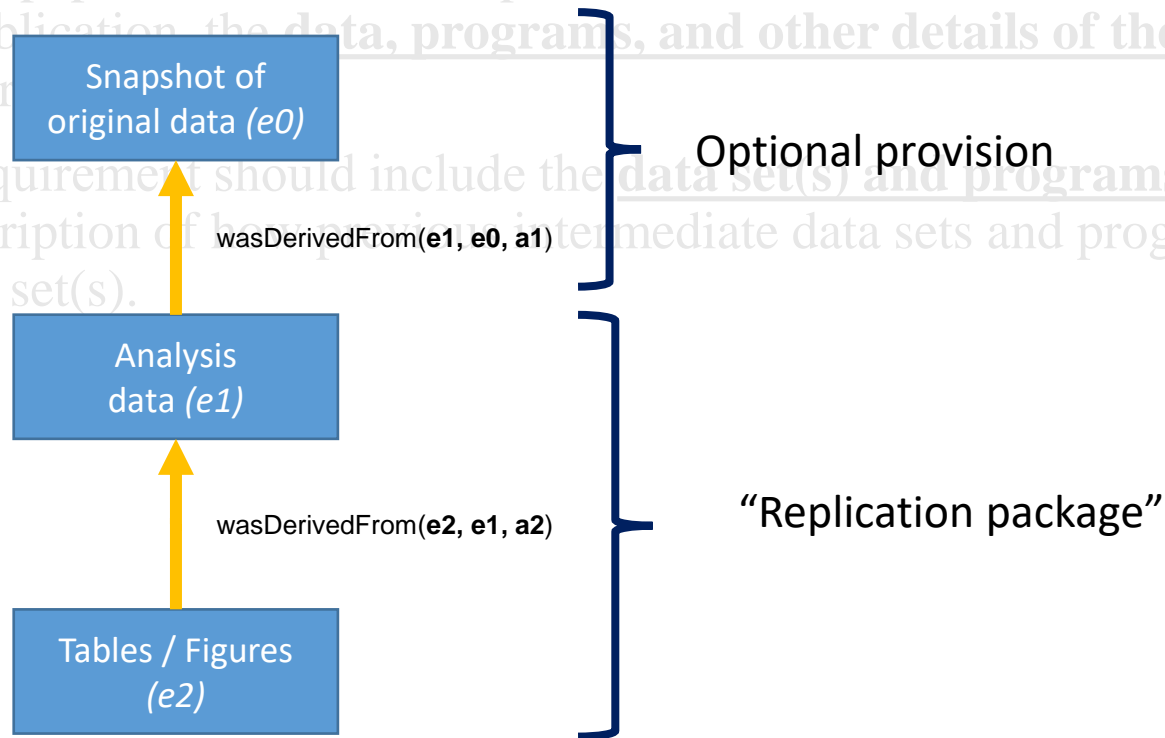
# AER policy ca. 2005

- It is the policy of the AEA to publish papers only if the **data used in the analysis** are clearly and precisely documented and are readily available to any researcher for purposes of replication.

- Authors of accepted papers that contain empirical work, simulations, or experimental work must provide, prior to publication, the **data, programs, and other details of the computations** sufficient to permit replication.

- … the minimum requirement should include the **data set(s) and programs used to run the final models**, plus a description of how previous intermediate data sets and programs were employed to create the final data set(s).
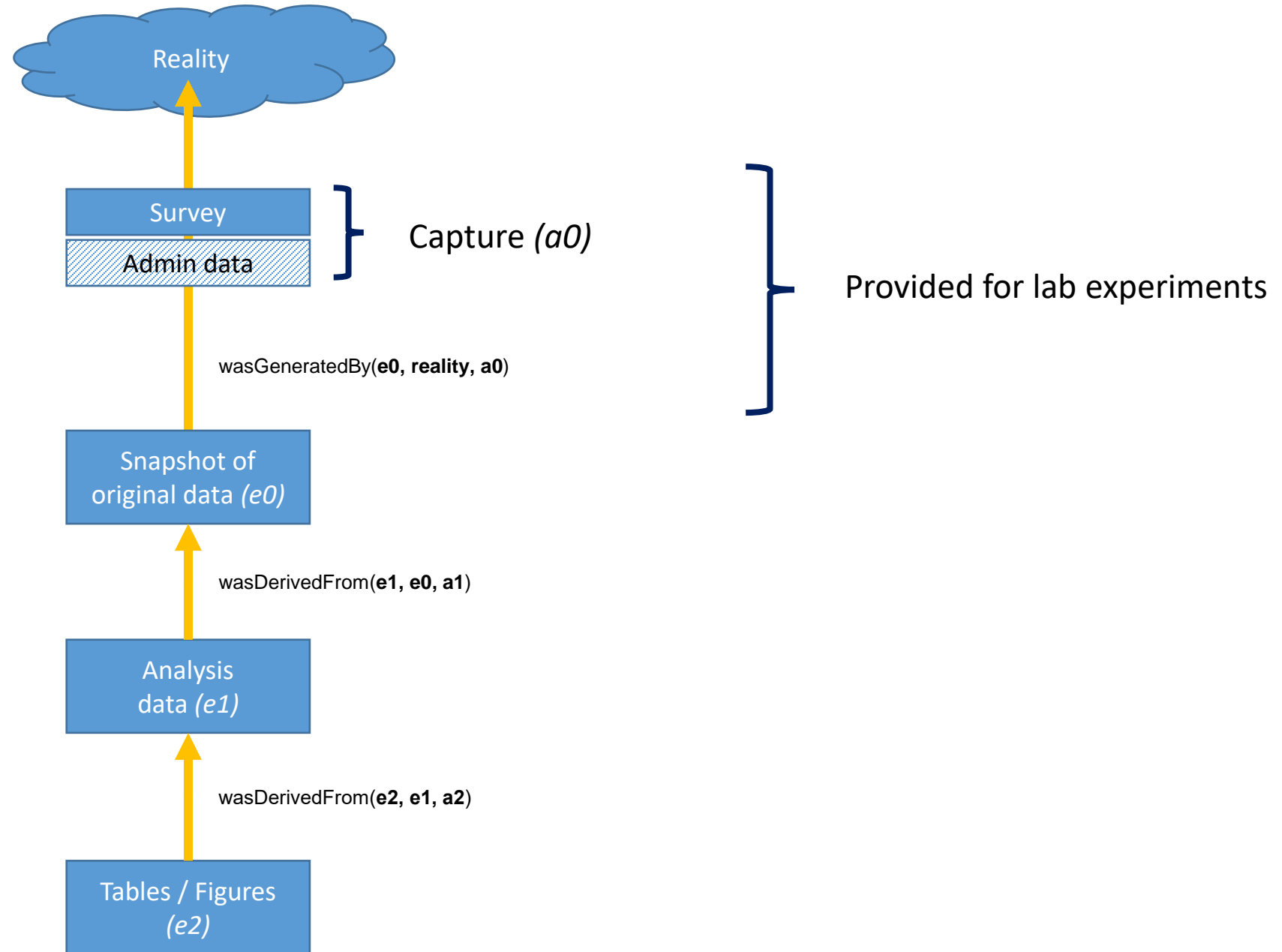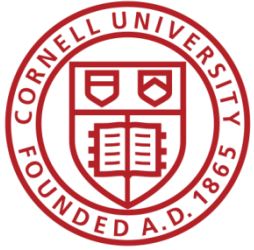


Analysis data *(e1)*

wasDerivedFrom(**e2, e1, a2**)

Tables / Figures *(e2)*

# AER policy ca. 2005

- It is the policy of the AEA to publish papers only if the **data used in the analysis** are clearly and precisely documented and are readily available to any researcher for purposes of replication.

- Authors of accepted papers that contain empirical work, simulations, or experimental work must provide, prior to publication, the **data, programs, and other details of the computations** sufficient to permit replication.

- … the minimum requirement should include the **data set(s) and programs used to run the final models**, plus a description of how previous intermediate data sets and programs were employed to create the final data set(s).

| Snapshot of original data *(e0)* |
| --- |

wasDerivedFrom(**e1, e0, a1**)

Optional provision

| Analysis data *(e1)* |
| --- |

wasDerivedFrom(**e2, e1, a2**)

"Replication package"

| Tables / Figures *(e2)* |
| --- |

### American Economic Review

The *American Economic Review* is a general-interest economics journal. Established in 1911, the *AER* is among the nation's oldest and most respected scholarly journals in economics.

### American Economic Review: Insights

*AER: Insights* is designed to be a top-tier, general-interest economics journal publishing papers of the same quality and importance as those in the *AER*, but devoted to publishing papers with important insights that can be conveyed succinctly.

### Journal of Economic Literature

The *Journal of Economic Literature* *(JEL)*, first published in 1969, is designed to help economists keep abreast of and synthesize the vast flow of literature.

### Journal of Economic Perspectives

The *Journal of Economic Perspectives* *(JEP)* fills the gap between the general interest press and academic economics journals.

### American Economic Journal: Applied Economics

*American Economic Journal: Applied Economics* publishes papers covering a range of topics in applied economics, with a focus on empirical microeconomic issues.

### American Economic Journal: Economic Policy

*American Economic Journal: Economic Policy* publishes papers covering a range of topics, the common theme being the role of economic policy in economic outcomes.

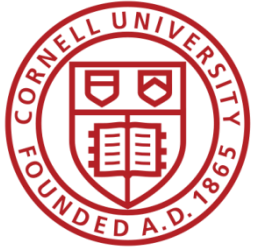### American Economic Journal: Macroeconomics

*American Economic Journal: Macroeconomics* focuses on studies of aggregate fluctuations and growth, and the role of policy in that context.

### American Economic Journal: Microeconomics

*American Economic Journal: Microeconomics* publishes papers focusing on microeconomic theory; industrial organization; and the microeconomic aspects of international trade, political economy, and finance.

# Replication continuum

| Same data | Same code | Same methods | Same context |
| --- | --- | --- | --- |
| | | | |
| | | | |

**Reproducibility**

**Replicability**

**Generalizability**

- Narrow Replication (Pesaran 2003)
- Pure Replication (Hamermesh 2007)
- Verification (Clemens 2015)

- Wide Replication (Pesaran 2003)
- Statistical Replication (Hamermesh 2007)
- Reproduction/Reanalysis (Clemens 2015)

- Wider Replication (Pesaran 2003)
- Scientific Replication (Hamermesh 2007)
- Reanalysis/Robustness (Clemens 2015)

# Historical challenges

# Poor citation practices

Provenance!

- **Macrodata:**

  "We use data downloaded from
  the Bureau of Economic Analysis…"

- **Microdata**:

  "… this paper uses data from
  the Current Population Survey…"

# Poor coding practices

- **Manual/non-automation**

  Code produces no meaningful output

- **Lack of robustness**:

  Bugs in the code

# Failure to curate

# AEA Data & Code Availability Policy (2019)

- It is the policy of the American Economic Association to publish papers only if the data used in the analysis are **clearly and precisely** documented and **access to the data and code is clearly and precisely documented and is non-exclusive to the authors.**

- Authors of accepted papers that contain empirical work, simulations, or experimental work must **provide**, **prior to acceptance**, the data, programs, and other details of the computations **sufficient to permit replication**, as well as **information about access to data and programs.**

# Current efforts at the AEA

- **Pre-emptively improve code archives**
  - By conducting reproducibility checks when we can
  - By working with groups that conduct reproducibility checks when we cannot
- **Better archives**
  - Greater transparency of the code and data archives
- **Better provenance tracking**
  - Leave code where it is when appropriate
  - Leave data where it is almost always
  - Display that information

# Current efforts at the AEA

- **Pre-emptively improve code archives**
  - By conducting reproducibility checks <small>when we can</small>
  - By working with groups that conduct reproducibility checks <small>when we cannot</small>
- **Better archives**
  - Greater transparency of the code and data archives
- **Better provenance tracking**
  - Leave code where it is when appropriate
  - Leave data where it is almost always
  - Display that information

perceived criteria of importance.

## 1. Importance

Data should be considered legitimate, citable products of research. Data should be accorded the same importance in the scholarly record as citat research objects, such as publications[1].

## 2. Credit and Attribution

Data citations should facilitate giving scholarly credit and normative and l attribution to all contributors to the data, recognizing that a single style or of attribution may not be applicable to all data[2].

## 3. Evidence

In scholarly literature, whenever and wherever a claim relies upon data, th corresponding data should be cited[3].

## 4. Unique Identification

A data citation should include a persistent method for identification that i actionable, globally unique, and widely used by a community[4].

## 5. Access

Data citations should facilitate access to the data themselves and to such

Data Citation Synthesis Group: Joint Declaration of Data Citation Principles. Martone M. (ed.) San Diego CA: FORCE11; 2014 [https://www.force11.org/group/joint-declaration-data-citation-principles-final].

# Remaining challenges

# Poor citation practices

- **Macrodata:**

  "We use data downloaded from
  the Bureau of Economic Analysis…"

- **Microdata**:

  "… this paper uses data from
  the Current Population Survey…"

# Poor coding practices

- **Manual/non-automation**

  Code produces no meaningful output

- **Lack of robustness**:

  Bugs in the code

# Failure to curate

# Challenges remaining

## Poor citation practices

- **Macrodata:**
  - "We use data downloaded from the Bureau of Economic Analysis…"
- **Microdata:**
  - "… this paper uses data from the Current Population Survey…"

## Poor coding practices

- **Manual/non-automation**
  - Code produces no meaningful output
- **Lack of robustness:**
  - Bugs in the code

## Failure to curate

- Required now
- Verified by Data Editor
- Not always in manuscript
- Challenging!

- Still an issue
- Self-check needed
- "Computational empathy" required

- Solved for code
- Fallback available for data
- Lots of issues with data without redistribution rights (confidential, otherwise)

# Actions
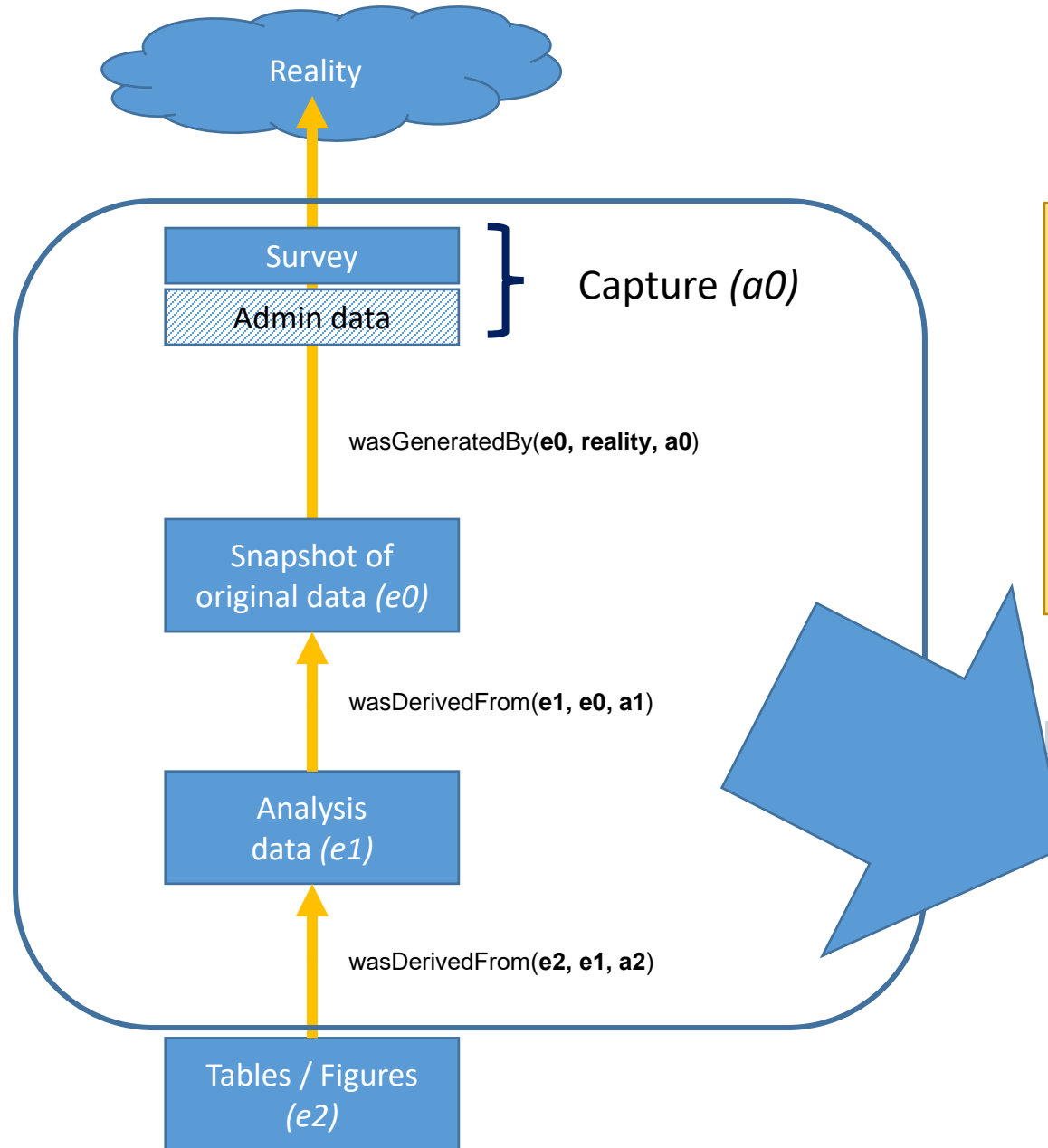
# Action: Reproducibility Check

## Data and Code Guidance by Data Editors

Guidance for authors wishing to create data and code supplements, and for replicators.

---

**Verification guidance**

On this page:
- Overview
- Review the README file
- For each listed data source
- For each listed table, figure, in-text number
- Conduct a code verification, if data is available
- Examples

**Overview**

This document describes

- what authors should check before providir
  journals
- what verifier teams should check for in the
  to them for the purpose of verification

# Stats on reproduced articles

Between Dec 1, 2019, and Nov 30, 2020, the AEA Data Editor team conducted

- **~832 assessments**

- ~**509 manuscripts** have been "accepted"

# Very little diversity in software

- **Stata** is the most popular statistical software in the journals of the AEA (**72.96%** of all supplements)

- followed by **Matlab** (**22.45%**)

# How this works

# Challenge: Where's the Data?

# Challenges remaining

### Poor citation practices

- **Macrodata:**
  "We use data downloaded from the Bureau of Economic Analysis…"
- **Microdata:**
  "… this paper uses data from the Current Population Survey…"

- Required now
- Verified by Data Editor
- Not always in manuscript
- Challenging!

### Poor coding practices

- **Manual/non-automation**
  Code produces no meaningful output
- **Lack of robustness:**
  Bugs in the code

- Still an issue
- Self-check needed
- "Computational empathy" required

### Failure to curate

Google
404. That's an error.
The requested URL /a_cool_website was not found on this server. That's all we know.

- Solved for code
- Fallback available for data
- Lots of issues with data without redistribution rights (confidential, otherwise)

# Three pieces

**Where did the author get the data?**

**Where can others get the same data?**

**Persistence of the data**

Reality

Survey

Admin data

Capture *(a0)*

wasGeneratedBy(**e0, reality, a0**)

Snapshot of
original data *(e0)*

wasDerivedFrom(**e1, e0, a1**)

Analysis
data *(e1)*

wasDerivedFrom(**e2, e1, a2**)

Tables / Figures
*(e2)*

Easy case:
Author-conducted survey

# Actions: Provide guidance

# Direct guidance to authors

# Enhanced guidance for authors



**Data and Code Guidance by Data Editors**

Guidance for authors wishing to create data and code supplements, and for replicators.

Cite this page as: Social Science Data Editors. 2021. "Guidance on Data Citations". *Data and Code Guidance by Data Editors.* Accessed at https://social-science-data-editors.github.io/guidance

## Guidance on Data Citations

On this page:
- What is not a data citation
  - Better
- Why data citations
- Generic Guidance
  - Websites
  - Online databases
  - Data distributed as supplementary data
- Specific Guidance
  - Producer
  - Distributor
  - Dates
  - Many related datasets
  - Offline access mechanism
  - Confidential databases
  - No formal access mechanism
  - Data provider cannot be named

One of the most vexing issues is how to cite data. This document goes through a few common scenarios not covered elsewhere.

### What is not a data citation



## A template README for social science replication packages.

The template README provided on this website is in a form that follows best practices as defined by a number of data editors at social science journals.

*Authors:* Lars Vilhuber, Miklos Kóren, Joan Llull, Marie Connolly, Peter Morrow

This project is maintained at social-science-data-editors/template_README

*Disclaimer*

DOI 10.5281/zenodo.4319999

## A template README for social science replication packages

The template README provided on this website is in a form that follows best practices as defined by a number of data editors at social science journals. A full list of endorsers is listed in Endorsers.

### Versions

The most recent version is available at https://social-science-data-editors.github.io/template_README/. Specific releases can be found at https://github.com/social-science-data-editors/template_README/releases.

### Formats

The template README is available in a variety of formats:

- HTML (best for reading)
- LaTeX
- Word
- PDF
- Markdown

### Description

The typical README in social science journals serves the purpose of guiding a reader through the available material and a route to replicating the results in the research paper, including the description of the origins of data and/or description of programs. As such, a good README file should first provide a brief overview of the available material and a brief guide as

# Edit and checks

In order to download the file you are asked to fill the following registration form and agree on the "Conditions of Use". Please read it carefully before proceeding to the download.

**PERSONAL DATA**

| | |
|---|---|
| Title (position): | |
| Full name: | |
| Company/Institution: | |
| E-mail: | |

**FILE USAGE**

| | |
|---|---|
| Project title: | |
| Intended use: | |

Brief description of the purpose of application:

**CONDITIONS OF USE**

1. Restrictions

These data files are available without restrictions, provided
a) that they are used for non-profit purposes; and
b) correct citations are provided and sent to the World Values Survey Association for each publication of results based in part or entirely on these data files. This citation will be made freely available; and
c) the data files themselves are not redistributed.

2. Correct citation

- **What does the site say?**

Please use the following citation when referring to this file in the different versions:
Inglehart, R., C. Haerpfer, A. Moreno, C. Welzel, K. Kizilova, J. Diez-Medrano, M. Lagos, P. Norris, E. Ponarin & B. Puranen et al. (eds.). 2014. World Values Survey: Round Six - Country-Pooled Datafile Version: www.worldvaluessurvey.org/WVSDocumentationWV6.jsp. Madrid: JD Systems Institute.
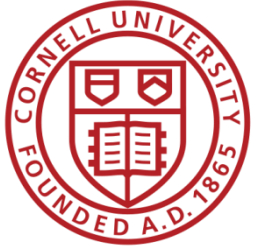
- **Is that in the README / Paper/ Appendix?**

- **Are all the conditions met/described?**

# Talking to Data Providers

- Federal statistical agencies
- Private data providers
- Individual research groups that produce data
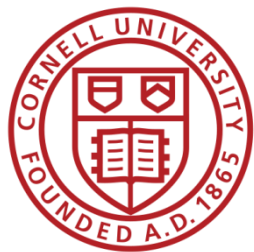
# The role for journals

# Data citations

Journals should uniformly
**monitor**
and
**enforce**
proper data citations!

# Goal: Transportability

Any standards, tools, methods: must be transportable across journals (no custom solutions)

# Social science "guild"



https://social-science-data-editors.github.io/guidance/

# Concluding remarks

wasDerivedFrom(e2, e1, a, g2, u1)

# Implicitly

- Goal to formally describe
  - Machine readable!
  - PIDs
  - Machine actionable?

# In practice

- Goal to formally describe
  - Machine readable!
  - PIDs
  - Machine actionable?

- Very few PIDs

- Data citations a challenge

- Click-through licenses

- Lack of licenses

- Lack of API

- Lack of persistent query extracts

- …

# Thank you