# Crowdsourcing Metadata – Challenges and Outlook
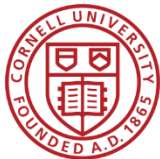
Washington, 10 May 2016

*Lars Vilhuber, William Block (Cornell University)*

# Crowdsourcing Metadata – Challenges and Outlook

Washington, 10 May 2016

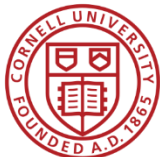*Lars Vilhuber, William Block (Cornell University)*

# Acknowledgements

Based on work with

- Benjamin Perry (formerly Cornell University)
- Venkata Kambhampaty (formerly Cornell University)
- Kyle Brumsted (McGill University)
- Jeremy Williams (Cornell University)
- Carl Lagoze (University of Michigan)
- John Abowd (Cornell University)

and materials presented in INFO 7470, all of that with funding by NSF Grant #1131848
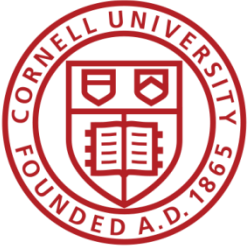
# I'm going to argue that…

- **Replicability** is a problem…
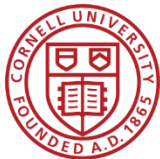  - and (A) easier **deposit** methods could alleviate it
  - but progress is slow

# I'm going to argue that…

- Having replicable archives **shifts** the problem…
  - in time: (B) older articles cannot be **linked to data**
  - in scope: (C) curators need **expert help** in documenting the data

# A test

LDI "reproducibility" project:
Kingi, Stanchi, Vilhuber (2016, unpublished)
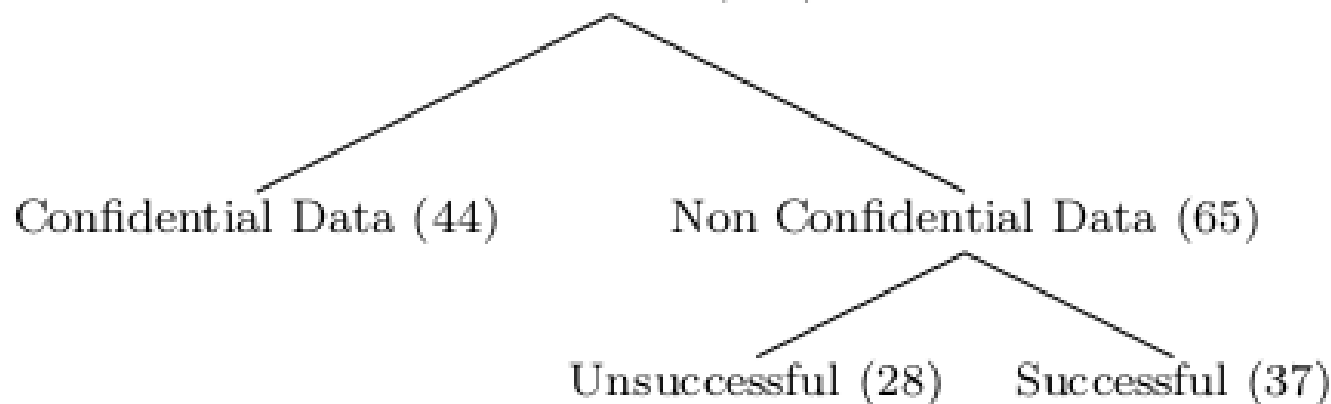"The Reproducibility of Economics Research"

# Kingi, Stanchi, Vilhuber (2016)

- 109 articles in American Economic Journal: Applied Economics
- Simpler test:

## **Do the provided data and programs yield the published results?**
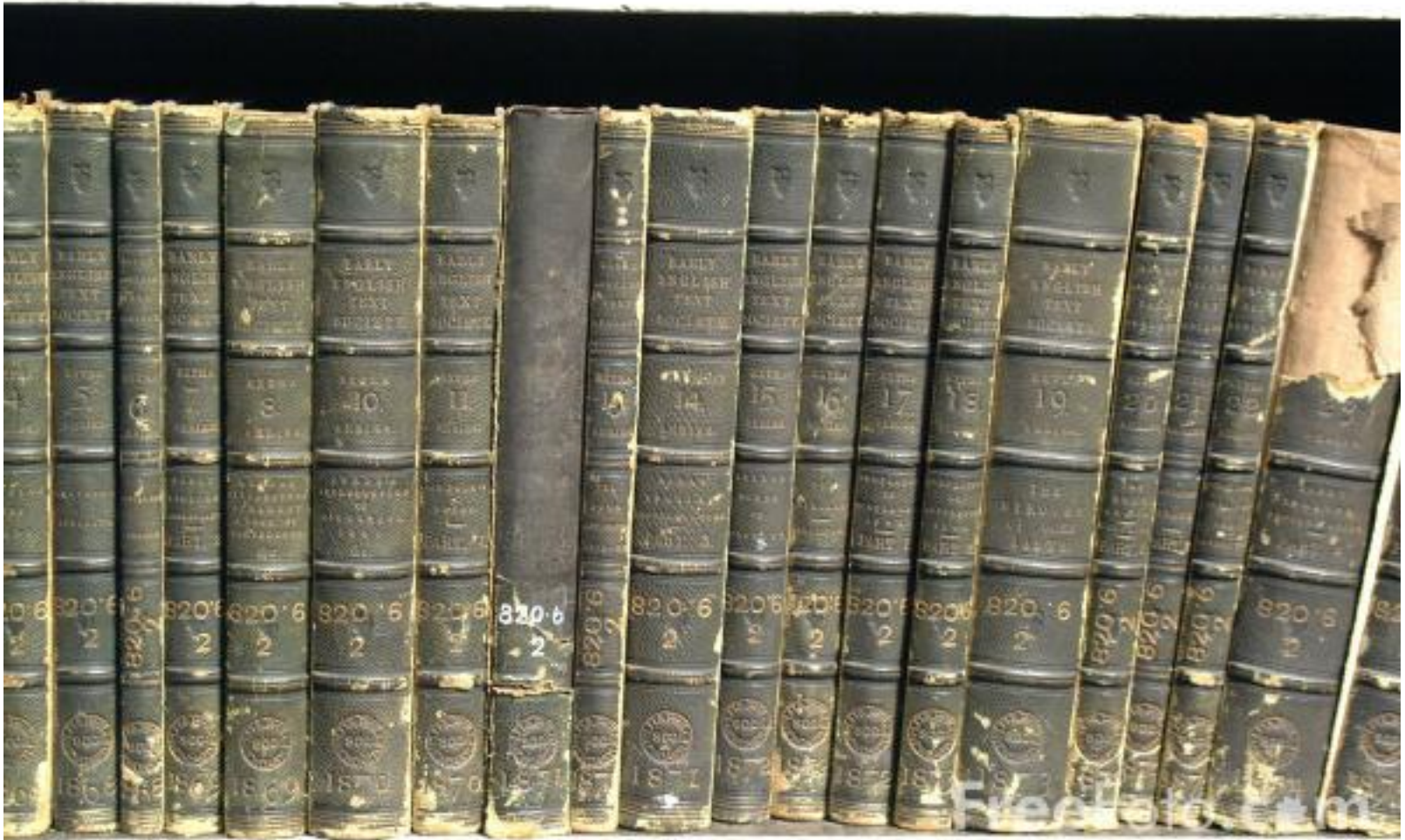
# Kingi, Stanchi, Vilhuber (2016)

Figure 1: A Breakdown of the Articles

Total Articles (109)

Confidential Data (44)   Non Confidential Data (65)

Unsuccessful (28)   Successful (37)

# The old source of knowledge – and data!

# Options are available

- Social and behavioral sciences

# Options are available

- "Research data"



The **Dataverse** Project

Open source research data repository software

**Researchers** — Enjoy full control over your data. Receive *web visibility, academic credit,* and *increased citation counts.* A personal dataverse is easy to set up, allows you to display your data on your personal website, can be branded uniquely as your research program, makes your data more discoverable to the research community, and satisfies data management plans. Want to set up your personal dataverse?

**Journals** — Seamlessly manage the submission, review, and publication of data associated with published articles. Establish an *unbreakable link* between *articles in your journal* and *associated data.* Participate in the open data movement by using Dataverse as part of your journal data policy or list of repository recommendations. Want to find out more about journal dataverses?
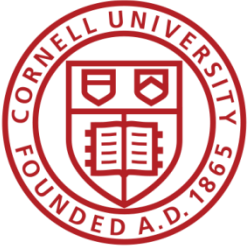
**Developers** — Participate in a vibrant and growing community that is helping to drive the norms for sharing, preserving, citing, exploring, and analyzing research data. Contribute code extensions, documentation, testing, and/or standards. *Integrate research analysis, visualization* and *exploration tools,* or other research and data archival systems with Dataverse. Want to contribute?

**Institutions** — Establish a research data management solution for your community. Federate with a growing list of Dataverse repositories worldwide for increased discoverability of your community's data. Participate in the drive to set norms for sharing, preserving, citing, exploring, and analyzing research data. Want to install a Dataverse repository?

# **Researchers don't use them**

Training? Incentives? Ease of use?

# A Good Example

# Gentzkow, Shapiro, Sinkinson (2014)

- Gentzkow, Matthew, Jesse M. Shapiro, and Michael Sinkinson. 2014. "*Competition and Ideological Diversity: Historical Evidence from US Newspapers.*"
  American Economic Review, 104(10): 3073-3114.
  DOI: 10.1257/aer.104.10.3073

- Data at http://doi.org/10.3886/E1361V3

Find data ▾   Share data ▾   Institutional repositories ▾

Browse Data  >  Circulation of US Daily Newspapers, 1924, Audit Bureau of Circulations.

# Circulation of US Daily Newspapers, 1924, Audit Bureau of Circulations.

Principal Investigator(s) :  Gentzkow, Matthew (University of Chicago. Booth School of Business); Shapiro, Jesse (University of Chicago. Booth

| Title | Date Entered | File Type |
|---|---|---|
| codebook<br>Gentzkow, Matthew; Shapiro, Jesse; Sinkinson, Michael | 2014-04-03<br>9:08 PM | .txt |
| data<br>Gentzkow, Matthew; Shapiro, Jesse; Sinkinson, Michael | 2014-04-03<br>9:09 PM | |
| Orig<br>Gentzkow, Matthew; Shapiro, Jesse; Sinkinson, Michael | 2014-04-07<br>1:29 PM | |

**Citation:** Gentzkow, Matthew; Shapiro, Jesse; Sinkinson, Michael. Circulation of US Daily Newspapers, 1924, Audit Bureau of Circulations.. Ann Arbor, MI: Inter-university Consortium for Political and Social Research [distributor], 2014-09-26. http://doi.org/10.3886/E1361V3

**Persistent URL:** http://doi.org/10.3886/E1361V3

## Project Description

### Summary

The focus of this data collection was the historical circulation and subscription prices of US daily newspapers in 1924. These data are obtained from audit reports obtained from the Audit Bureau of Circulations, an independent organization created to verify circulation. They include circulation by town and delivery channel for each newspaper.
The sample is all audited daily newspapers by the Audit Bureau of Circulations.
All pdfs and extracted .dtas. Copyright belongs to the Audit Bureau of Circulations. We have obtained written permission from the Audit Bureau of Circulations to post the PDFs and data files.

# What's good about this?

- Permanent URL
- Availability of
  - Original data
  - Transformed data
  - Open availability
- Easy online inspection

# Not Perfect

- Archive at openICPSR not actually tied to article and vice-versa

and Exit on Electoral Politics." *American Economic Review* 101 (7): 2980–3018.

**Gentzkow, Matthew, Jesse M. Shapiro, and Michael Sinkinson.** 2014. "Competition and Ideological Diversity: Historical Evidence from US Newspapers: Dataset." *American Economic Review.* http://dx.doi.org/10.1257/aer.104.10.3073.

▶**George, Lisa, and Joel Waldfogel.** 2003. "Who Affects Whom in Daily Newspaper Markets?" *Journal of Political Economy* 111 (4): 765–84

- Conversely, "online appendix" just a "blob"

Article Full-Text Access

Full-text Article

**Additional Materials**

Download Data Set (1.99 GB) | Online Appendix (148.69 KB) | Author Disclosure Statement (10.33 KB)

**Authors**

# Journals are starting to use them

# But the biggest problem…

Figure 1: A Breakdown of the Articles

Total Articles (109)

Confidential Data (44)

Non Confidential Data (65)

Unsuccessful (28)    Successful (37)

# Back to confidential data

- Articles using confidential data are (weakly) more cited than others

# But: for confidential data…

- Data is not available

- Metadata is not available

- Programs? So-so…

# Should We Just Trust These Guys?

# Some are quite commendable

# Even detailed information

**Full Title**
General Social Survey, 2011: Cycle 25, Family

**Subtitle**
Cycle 25, Family

**Alternative Title**
GSS 2011: Family

**Parallel Title**
Enquête sociale générale, 2011: Cycle 25, Famille

**Identification Number**
ca-statcan-68196

**Authoring Entity**

| Name | Affiliation |
|---|---|
| Statistics Canada | StatCan |

**Producer**

| Name | Affiliation | Abbreviation | Role |
|---|---|---|---|
| Statistics Canada | | StatCan | |

**Copyright**
Copyright © Statistics Canada, 2012

**Date of Distribution**
2012-07-18

**Series Information**
General Social Survey - Family (GSS) [4501]

**Version**
16769.6

# How many users actually use that?

# Data documentation is dry

- How reliable is that question?

Dataset: General Social Survey, 2011: Cycle 25, Family

Cycle 25, Family

## Variable PA_Q240: Year parents separated

**LITERAL QUESTION**
In what year did your parents separate?

**Values  Categories**
9997     Not asked
9998     Not stated
9999     Don't know

**SUMMARY STATISTICS**
This variable is numeric

**UNIVERSE**
Respondents who answered: PA_Q230 = 1.

**NOTES**
This variable is suppressed on the public use microdata file.

# Don't (just) liberate the data!

Liberate the data users!

# Our contribution

Leverage researcher knowledge

# Our Approach

- Rely on open standards, namely the Data Documentation Initiative (DDI) schema

- Provide easy-to-use tools and interfaces to structured metadata

- Build infrastructure that enables data curators to leverage community-driven input to official documentation

31

# How?

# CED²AR

The Comprehensive Extensible Data Documentation and Access Repository

# What is CED²AR?

- Metadata curation software
- Designed for documenting existing datasets
- Funded by NSF grant #1131848
- Online at www2.ncrn.cornell.edu/ced2ar-web

# What is CED²AR?

# Basic Information Flow



*Staging Area*

*Public Facing*

Datasets → Internal Metadata

Official Metadata

Crowdsourced Metadata

*User switches*

# Basic Information Flow

*Staging Area*

*Public Facing*

Datasets → Internal Metadata

Official Metadata

*User switches*

Crowdsourced Metadata

# Internal Processing

1. Creation of skeletal metadata
   - Assuming data is already curated
   - Converting <u>data</u> into standardized <u>metadata</u>
     - Tools included (for SAS, Stata, SPSS, CSV), not discussed here

2. Hand editing and subsetting
   - Adding verbose descriptions
   - Applying disclosure limitation

3. User accessible
   - These tools can be manipulated by normal users
   - They could be incorporated into existing workflows

# Internal Processing

- Simple editing interface
  - Web-based, with limited rich text features
  - Math allowed (LaTeX)
- Feedback
  - Completeness of codebook?
  - Without technical jargon!
  - Can be tuned

# Internal Processing: Hand Editing

## Abstract

The SIPP Synthetic Beta (SSB) is a Census Bureau product that integrates person-level micro-data from a household survey with administrative tax and benefit data. These data link respondents from the Survey of Income and Program Participation (SIPP) to Social Security Administration (SSA)/Internal Revenue Service (IRS) Form W-2 records and SSA records of retirement and disability benefit receipt, and were produced by Census Bureau staff economists and statisticians in collaboration with researchers at Cornell University, the SSA and the IRS. The purpose of the SSB is to provide access to linked data that are usually not publicly available due to confidentiality concerns.

To overcome these concerns, Census has synthesized, or modeled, all the variables in a way that changes the record of each individual in a manner designed to preserve the underlying covariate relationships between the variables. The only variables that were not altered by the synthesis process and still contain their original values are gender and a link to the first reported marital partner in the survey. Seven SIPP panels (1990, 1991, 1992, 1993, 1996, 2001, 2004) form the basis for the SSB, with a large subset of variables available across all the panels selected for inclusion and harmonization across the years. Administrative data were added and some editing was done to correct for logical inconsistencies in the IRS/SSA earnings and benefits data.

p

*This field supports ASCII math See FAQ for details.*

# Internal Processing: Scoring

- Provide feedback to improve sparse documentation



CED2AR / SIPP Synthetic Beta v6 / Score

## Codebook Score

### Variables

100.0% of variables have labels

85.1% of variables have significant full descriptions
*Variables without significant full descriptions ... more*

43.0% of variables have values
*Variables without values ... more*

0.0% of variables have summary statistics

### Title Page

Missing related studies
Missing access conditions
Missing bibliographic citation
Missing related publications

## Overall Score

80.3%

# Fine-grained access controls

Important when working with confidential (meta)data

# Internal Processing: Access Control

- Marking elements with different restrictions

Select what sub-elements to mark
☐ Select All
  ☐ Mean                              ☐ Values
  ☐ Median                            ☑ Value Frequencies
  ☐ Mode                               ☑ Value Percentages
  ☐ Valid                               ☑ Value Crosstabs
  ☐ Invalid                            ☑ Other Value Statistics
  ☑ Min                                ☐ Label
  ☑ Max                              ☐ Notes
  ☐ Standard Deviation
  ☐ Other Summary Statistics

Select what access level to apply, then check which variables to apply to. Finally, click changes levels.

| restricted ▼ | Change Levels |

| | Variable Name | Label | Top Access Level |
|---|---|---|---|
| ☑ | afdc_MN | Indicator for receipt of AFDC or TANF benefits | released |
| ☑ | afdcamt_MN | Amount of AFDC received | released |
| ☐ | birthdate | Date of Birth | released |
| ☐ | current_enroll_coll | Currently Enrolled in College | released |
| ☐ | current_enroll_hs | Currently Enrolled in HS (or less) | released |

42

# Workflow control

- Ability to view additions/subtractions
    - Between versions
    - Between crowd-sourced information and official information
- Ability to control access
    - Editing versus viewing
    - Authentication and reputation

# Versioning

## All changes are logged externally via Git

### Commits

| Author | Commit | Message | | Date |
|--------|--------|---------|------|------|
| ? tomcat7 | 0fea515 | {ssbv601,lars@vilhuber.com,cover} | ⑂ ssbtesting | 37 minutes ago |
| ? tomcat7 | 5e824de | {ssbv601,lars@vilhuber.com,cover}{ssbv601,lars@vilhuber.com,var,fl... | ⑂ ssbtesting | an hour ago |
| ? tomcat7 | c03c50f | Commiting codebooks retrieved directly from BaseX | ⑂ cestesting | 4 days ago |
| ? venkata | a61abe3 | {testlbdv1,anonymous,edit} | ⑂ vrk4 | 4 days ago |
| ? venkata | 5b1e51e | {testlbdv1,anonymous,edit} | ⑂ vrk4 | 4 days ago |
| ? tomcat7 | 5edbff9 | {acs2009,bap63@cornell.edu,edit} | ⑂ cestesting | 5 days ago |
| ? tomcat7 | d66d3d4 | {ssbv601,lorireeder@gmail.com,var,phus_ssdi_benefit_totamt_k}{ss... | ⑂ ssbtesting | 5 days ago |
| ? tomcat7 | 1f845c1 | {siabv1,warren.brown48@gmail.com,cover}{siabv1,bap63@cornell.e... | ⑂ cestesting | 5 days ago |
| ? tomcat7 | eb77f31 | {siabv1,warren.brown48@gmail.com,var,bild}{siabv1,warren.brown4... | ⑂ cestesting | 2015-11-17 |
| ? tomcat7 | b34a118 | {siabv1,warren.brown48@gmail.com,var,persnr}{siabv1,warren.brow... | ⑂ cestesting | 2015-11-17 |
| ? venkata | 2cb6d7d | {lbdv2,anonymous,cover} | ⑂ vrk4 | 2015-11-17 |
| ? tomcat7 | 1263bcf | {siabv1,warren.brown48@gmail.com,var,bnn}{siabv1,warren.brown4... | ⑂ cestesting | 2015-11-17 |
| ? tomcat7 | aaf94f1 | {siabv1,bap63@cornell.edu,edit}{blss2011,bap63@cornell.edu,edit}{... | ⑂ cestesting | 2015-11-17 |
| ? tomcat7 | 0e52a6e | Commiting codebooks retrieved directly from BaseX | ⑂ cestesting | 2015-11-17 |

# Basic Information Flow

*Staging Area*

*Public Facing*



Datasets → Internal Metadata

Official Metadata

*User switches*

Crowdsourced Metadata

# Basic Information Flow

*Staging Area*

*Public Facing*

Datasets → Internal Metadata

Official Metadata

*User switches*

Crowdsourced Metadata

# Official view

## CED²AR

Official Server - The Comprehensive Extensible Data Documentation and Access Repository

Search Variables    Browse Variables ▾

You are viewing the official *crowdsourced contributions.*

CED2AR / SIPP Synthetic Beta

# SIPP Synthetic Beta
## v6.02

View Variables *(123 variables)*
Last update to metadata: 2015-11-24 10:05:15 (upload date)
Document Date: November 12, 2015

Codebook prepared by: Cornell NSF-Census Research Network

Data prepared by: United States Department of Commerce. Bureau of the Census.

## Data Distributed by:

Labor Dynamics Institute
http://www2.vrdc.cornell.edu/news/data/sipp-synthetic-beta-file/

United States Department of Commerce. Bureau of the Census.
http://www.census.gov/programs-surveys/sipp/methodology/sipp-synthetic-beta-data-product.html

49

# Crowdsourced view

## CED²AR

Community Development Server (Beta) - The Comprehensive Extensible Data
Documentation and Access Repository

*You are viewing crowdsourced metadata. View the official version .*

## SIPP Synthetic Beta

→] 🖶 </> SAS Stata 📄

### v6.02

View Variables *(123 variables)*
Last update to metadata: 2015-11-24 09:59:07 (auto-generated)
Document Date: November 12, 2015

Codebook prepared by: Cornell NSF-Census Research Network

Data prepared by: United States Department of Commerce. Bureau of the Census.

### Data Distributed by:

Labor Dynamics Institute
http://www2.vrdc.cornell.edu/news/data/sipp-synthetic-beta-file/

United States Department of Commerce. Bureau of the Census.
http://www.census.gov/programs-surveys/sipp/methodology/sipp-synthetic-beta-data-product.html

50

# Authentication and Attribution

- When opening up contributions to a wide audience, how to triage between "rants" and meaningful contributions?

- Here: Use of ORCID (academic network) for authentication

- Public attribution with link to (verified) academic ID is key for positive feedback (your effort is recognized) and prevention of negative contribution (your rant is traceable to you!)

# Authentication

- Supports OpenID and OAuth2
  - Currently using Google and ORCID with OAuth2
  - Developing connectors to work with additional providers
- CED$^2$AR handles identity management

**Login to Continue**

Please choose authentication method

g+ Google     ORCID

# Editing made easy

You are viewing crowdsourced metadata. View the *official version* .

CED2AR / SIPP Synthetic Beta v6.02

# SIPP Synthetic Beta v6.02 ✎ ⓘ

View Variables *(123 variables)*
View codebook score
Last update to metadata: 2016-01-26 14:36:26 (auto-generated) ⓘ

Document Date: November 12, 2015 ✎ ⓘ

Codebook prepared by: Cornell NSF-Census Research Network ✎ ✚ ⓘ

Data prepared by: United States Department of Commerce. Bureau of the Census. ✎ ✚ ⓘ

## Data Distributed by: ⓘ

Labor Dynamics Institute ✎ 🗑

http://www2.vrdc.cornell.edu/news/data/sipp-synthetic-beta-file/ ✎

United States Department of Commerce. Bureau of the Census. ✎ 🗑

http://www.census.gov/programs-surveys/sipp/methodology/sipp-synthetic-beta-data-product.html ✎
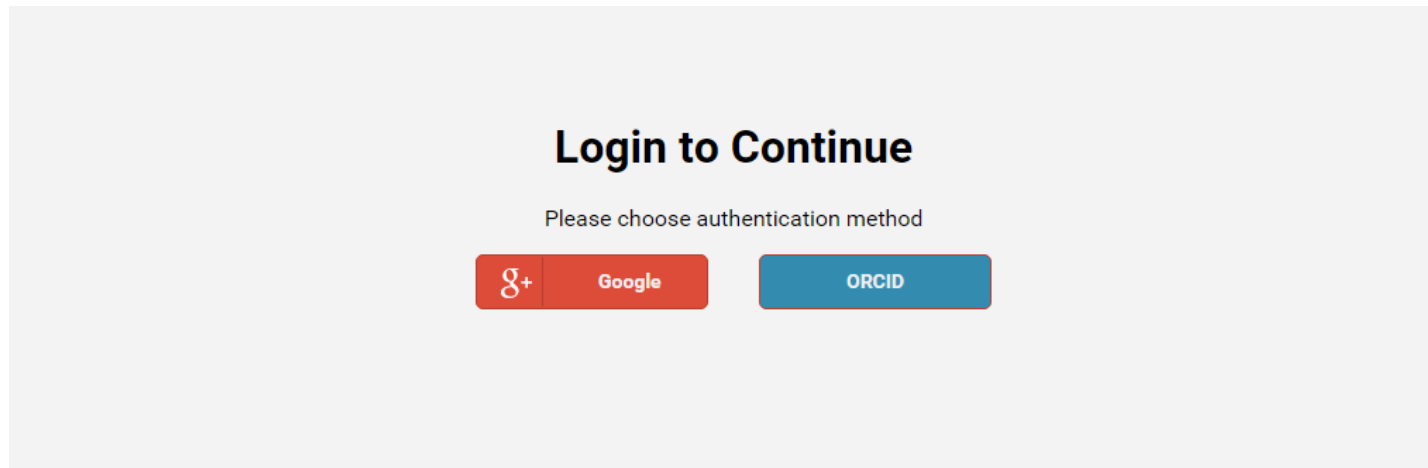
CED2AR / SIPP Synthetic Beta v6.02 / totearn_ser_YYYY

## Variable Name ❶

totearn_ser_YYYY

## Top Access Level

released ✎

## Label

SER: Capped Earnings from all FICA-covered jobs ✎

Access Level: released ✎

## Codebook

SIPP Synthetic Beta v6.02

## Concept ❶

✎

## Type

numeric

## Question Text ➕ ❶

## Full Description ✎ ❶

Person-level annual earnings that were taxed by FICA; these variables include earnings only up to the FICA taxable maximum and cover the years 1951-2011. These earnings are the inputs for calculating the OASDI benefit a person and his or her spouse will receive upon retirement or disability.

## Files ❶

ssb_v6_0_synthetic1_1.sas7bdat  http://www.census.gov/programs-surveys/sipp/methodology/sipp-synthetic-beta-data-product.html ☐ ( SAS )

ssb_v6_0_synthetic1_1.dta  http://www.census.gov/programs-surveys/sipp/methodology/sipp-synthetic-beta-data-product.html ☐ ( Stata )

# Notes ⓘ

Note #1 - Access Level:  released   ✎

Having a 0 value on totearn_ser_yyyy could mean a couple of things: 1) this individual had no FICA-covered earnings in that year; 2)this individual had no labor income at all in this tax year; 3) this individual worked for an employer that failed to report earnings in this year (that is to say, this has nothing to do with whether a person filed taxes because the earnings are reported by the employer, not the employee). Prior to 1978, if a person has $0 earnings on the Summary Earnings Record (SER), there's really no way of knowing whether they had no earnings or whether they had non-FICA earnings because the SER only reports FICA-covered earnings reported by employers. For years 1978 and later, you can compare the SER to the Detailed Earnings Record (DER). The DER captures all earnings subject to income tax, so both FICA and non-FICA earnings are reported on the DER.

If you are looking at earnings in earlier years, particular the 1960s and earlier, there will be more people with $0 earnings because many jobs were not FICA-taxable then. Even today, there are some instances of legitimate non-FICA earnings that would not be reflected on the SER. One example of this is that graduate student stipends are not taxed for FICA or Medicare, so these earnings would not be reflected on the SER (https://www.irs.gov/Charities--Non-

Profits/Student-Exception-to-FICA-Tax).   ✎

➕  Add Note

Type          numeric

# Notes

Save

Having a 0 value on totearn_ser_yyyy could mean a couple of things: 1) this individual had no FICA-covered earnings in that year; 2)this individual had no labor income at all in this tax year; 3) this individual worked for an employer that failed to report earnings in this year (that is to say, this has nothing to do with whether a person filed taxes because the earnings are reported by the employer, not the employee). Prior to 1978, if a person has $0 earnings on the Summary Earnings Record (SER), there's really no way of knowing whether they had no earnings or whether they had non-FICA earnings because the SER only reports FICA-covered earnings reported by employers. For years 1978 and later, you can compare the SER to the Detailed Earnings Record (DER). The DER captures all earnings subject to income tax, so both FICA and non-FICA earnings are reported on the DER.

If you are looking at earnings in earlier years, particular the 1960s and earlier, there will be more people with $0 earnings because many jobs were not FICA-taxable then. Even today, there are some instances of legitimate non-FICA earnings that would not be reflected on the SER. One example of this is that graduate student stipends are not taxed for FICA or Medicare, so these earnings would not be reflected on the SER (https://www.irs.gov/Charities--Non-Profits/Student-Exception-to-FICA-Tax).

p

*This field supports ASCII math See* FAQ *for details.*

# Basic Information Flow

*Staging Area*

*Public Facing*



Datasets → Internal Metadata → Official Metadata

Official Metadata ⇅ Crowdsourced Metadata

*User switches*

# Everybody can see changes

## Remote

| | |
|---|---|
| Variable Name | totearn_ser_YYYY |
| Label | SER: Capped Earnings from all FICA-covered jobs |
| Codebook | SIPP Synthetic Beta v6.02 |
| Concept | |
| Concept Vocabulary | |
| Concept Vocabulary URI | |
| Type | numeric |
| Files | |

ssb_v6_0_2_syntheticK_M.sas7bdat
http://www.census.gov/programs-surveys/sipp/methodology/sipp-synthetic-beta-data-product.html
SAS

ssb_v6_0_2_syntheticK_M.dta
http://www.census.gov/programs-surveys/sipp/methodology/sipp-synthetic-beta-data-product.html
Stata

### Question Text

### Full Description

Person-level annual earnings that were taxed by FICA; these variables include earnings only up to the FICA taxable maximum and cover the years 1951-2011. These earnings are the inputs for calculating the OASDI benefit a person and his or her spouse will receive upon retirement or disability.

### Notes (0 total)

## Current

| | |
|---|---|
| Variable Name | totearn_ser_YYYY |
| Label | SER: Capped Earnings from all FICA-covered jobs |
| Codebook | SIPP Synthetic Beta v6.02 |
| Concept | |
| Concept Vocabulary | |
| Concept Vocabulary URI | |
| Type | numeric |
| Files | |

ssb_v6_0_synthetic1_1.sas7bdat
http://www.census.gov/programs-surveys/sipp/methodology/sipp-synthetic-beta-data-product.html
SAS

ssb_v6_0_synthetic1_1.dta
http://www.census.gov/programs-surveys/sipp/methodology/sipp-synthetic-beta-data-product.html
Stata

### Question Text

### Full Description

Person-level annual earnings that were taxed by FICA; these variables include earnings only up to the FICA taxable maximum and cover the years 1951-2011. These earnings are the inputs for calculating the OASDI benefit a person and his or her spouse will receive upon retirement or disability.

### Notes (1 total)

#1

Having a 0 value on totearn_ser_yyyy could mean a couple of things: 1) this individual had no FICA-covered earnings in that year; 2)this individual had no labor income at all in this tax year; 3) this individual worked for an employer that failed to report earnings in this year (that is to say, this has nothing to do with whether a person filed taxes because the earnings are reported by the employer, not the employee). Prior to 1978, if a person has $0 earnings on the Summary Earnings Record (SER), there's really no way of knowing whether they had no earnings or whether they had non-FICA earnings because the SER only reports FICA-covered earnings reported by employers. For years 1978 and later, you can

# Combining Knowledge: Merging

- Curators are given an interface to merge crowdsourced documentation with official

## Merge Variables

The following variables have changed:

cur_endmar

birthdate

Continue

# Combining Knowledge: Merging

## current_enroll_coll

### Crowdsourced Documentation

| | |
|---|---|
| Variable Name | current_enroll_coll |
| Label | Currently Enrolled in College |
| Codebook | SIPP Synthetic Beta v6 |
| Concept | |
| Concept Vocabulary | |
| Concept Vocabulary URI | |
| Type | numeric |
| Files | |

### Official Documentation

| | |
|---|---|
| Variable Name | current_enroll_coll |
| Label | ☐ Use crowdsourced  ☐ Use original  Currently Enrolled |
| Codebook | SIPP Synthetic Beta v6 |
| Concept | |
| Concept Vocabulary | |
| Concept Vocabulary URI | |
| Type | numeric |
| Files | |

# Combining Knowledge: Merging

## Crowdsourced Documentation

Last update to metadata: 2015-08-18 08:43:01 (upload date)

Document Date:

June 15**8**, 2014

### Citation

*Please cite this codebook as:*

Comprehensive Extensible Data Documentation and Access Repository. Codebook for the SIPP Synthetic Beta 5.1 [Codebook file]. Cornell Institute for Social and Economic Research and Labor Dynamics Institute [distributor]. Cornell University, Ithaca, NY, 2013

*Please cite this dataset as:*

U.S. Census Bureau. SIPP Synthetic Beta: Version 5.1 [Computer file]. Washington DC; Cornell University, Synthetic Data Server [distributor], Ithaca, NY, 2013

### Abstract

The SIPP Synthetic Beta (SSB) is a Census Bureau product that integrates person-level micro-data from a household survey with administrative tax and benefit data. These data link respondents from the Survey of Income and Program Participation (SIPP) to Social Security Administration (SSA)/Internal Revenue Service (IRS) Form W-2 records and SSA

## Official Documentation

Last update to metadata: 2015-10-23 11:12:44 (auto-generated)

Document Date:

☐ Use crowdsourced    ☐ Use original    June 15, 2014

### Citation

*Please cite this codebook as:*

Comprehensive Extensible Data Documentation and Access Repository. Codebook for the SIPP Synthetic Beta 5.1 [Codebook file]. Cornell Institute for Social and Economic Research and Labor Dynamics Institute [distributor]. Cornell University, Ithaca, NY, 2013

*Please cite this dataset as:*

U.S. Census Bureau. SIPP Synthetic Beta: Version 5.1 [Computer file]. Washington DC; Cornell University, Synthetic Data Server [distributor], Ithaca, NY, 2013

### Abstract

☐ Use crowdsourced    ☐ Use original

The SIPP Synthetic Beta (SSB) is a Census Bureau product that integrates person-level micro-data from a household survey with administrative tax and benefit data. These data link respondents from the Survey of Income and Program Participation (SIPP) to Social

62

# Combining Knowledge: Citations

- Contributors can be tracked for each of their changes

# Combining Knowledge: Citations

ORCID

1,757,580 ORCID iDs and counting. See more...

**Lars Vilhuber**

**ORCID ID**

iD orcid.org/0000-0001-5733-8932

> Education (3)

> Employment (1)

> Funding (7)

✔ Works (29)

↓↑ Sort

CED²AR: The Comprehensive Extensible Data Documentation and Access Repository

IEEE/ACM Joint Conference on Digital Libraries

2014-09 | conference-paper

DOI: 10.1109/jcdl.2014.6970178

Source: CrossRef Metadata Search                    ✔ Preferred source
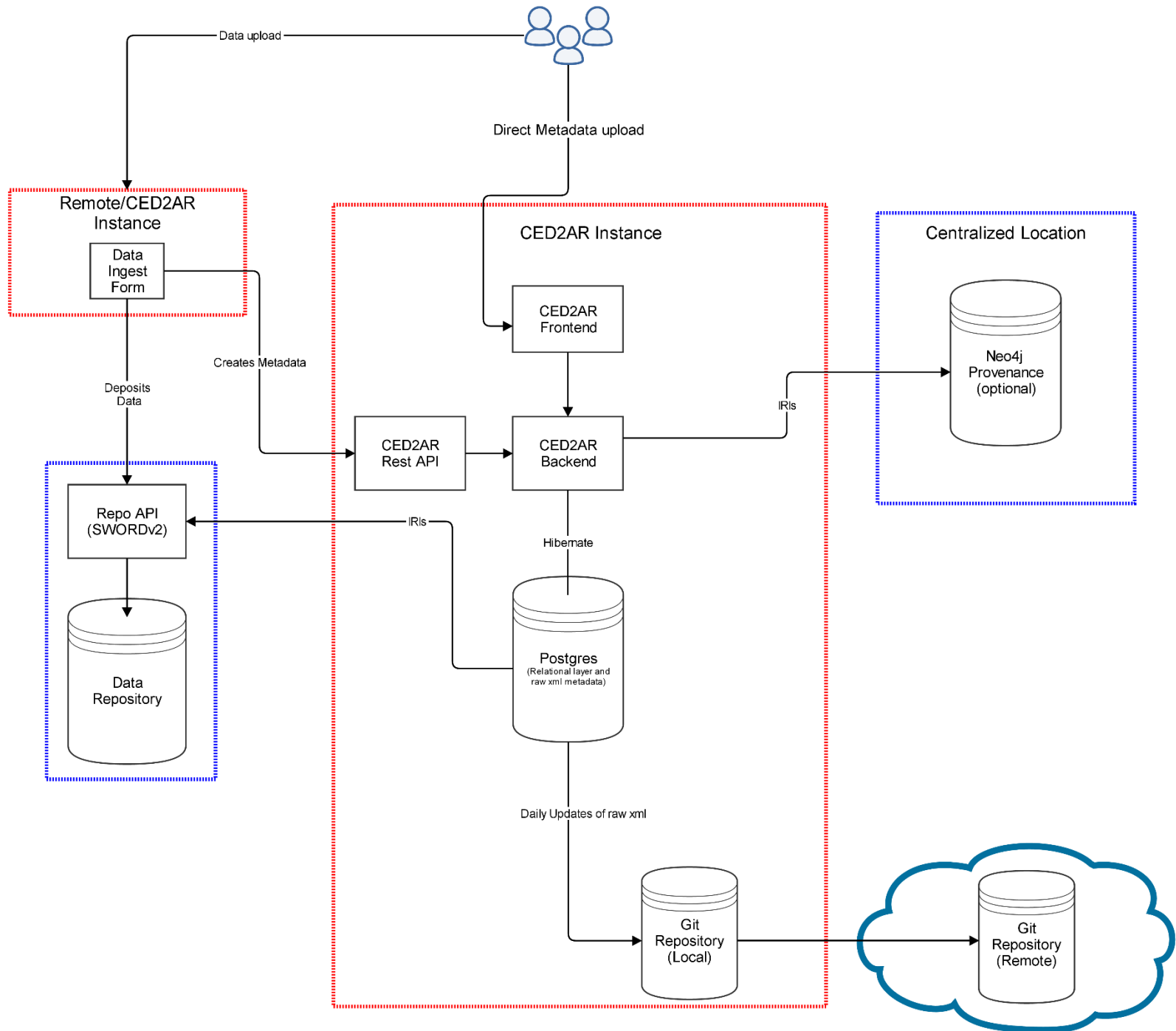
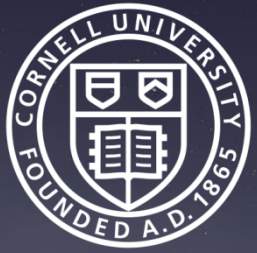# Next steps

# Making CED²AR v2 robust

- Addition of UTF-8 editing support (Spanish, French, Portuguese, etc.)

- Additional fields (link to survey questions, anything within DDI-C)

- Bug fixes

# Making CED²AR scalable in V3

- Current implementation of CED²AR is packaged for a single server (=portable)
  - Already a scalable archive backend (git)
  - Could be fronted by a load balancer
  - But live server is not scalable
- Current implementation of CED²AR is tied to a limited schema – adding schema is hard

Thank you!
Questions?

ced2ar-devs-l@cornell.edu