

CISER

Cornell Institute for Social and Economic Research
A LEADER IN SOCIAL SCIENCE DATA AND COMPUTING



The Complicated Provenance of American Community Survey Data: How Far will PROV and DDI Take Us?

William C. Block,¹ Warren Brown,¹ Jeremy Williams,¹ Lars Vilhuber,² and Carl Lagoze,³

¹ Cornell Institute Social and Economic Research (CISER), Cornell University

² Labor Dynamics Institute (LDI), Cornell University

³ School of Information, University of Michigan

— Presentation at the 2nd Annual North American DDI User Conference (NADDI14)
Vancouver, British Columbia, Canada
2 April 2014



Cornell University

Outline

Answering the Q: How far will PROV and DDI take us? We don't know; complicated story!

- Background/Previous Work
- Use Case(s) involving ANCESTRY Variable in ACS
- Technical solutions at File (Dataset) and Variable Level
- Future Work

Questions and Discussion

CISER

NSF-Census Research Network (NCRN) – Cornell Node
 (“Integrated Research Support, Training and
 Documentation”)

- CED²AR is one part of this project
- Funded by NSF Grant [#1131848](#).
- For more information, see www.ncrn.cornell.edu.



CISER

Part of NCRN Research Network



NSF-Census Research Network

[Home](#)

[News](#)

[Events](#)

[Documents](#)

[Nodes](#)

[About](#)

[Contact](#)



Eight nodes comprised of researchers conducting innovative, high-disciplinary investigations of theory, methodology and computational tools of interest.

[NCRN Coordinating Office](#)

[Carnegie-Mellon University](#)

[Cornell University](#)

[Duke University / National Institute of Statistical Sciences \(NISS\)](#)

[Northwestern University](#)

[University of Colorado at Boulder / University of Tennessee](#)

[University of Michigan](#)

[University of Missouri](#)

[University of Nebraska](#)

Latest News

[University of Colorado at Boulder Sharing Video about NCRN Research](#)

MARCH 19, 2014

[Registration Now Open for NCRN Spring Workshop](#)

MARCH 12, 2014

Calendar of events

APR
2

3:00pm to
4:30pm

[NCRN Virtual Seminar - Survey Informatics: Ideas, Issues, and Opportunities](#)

Speakers: [Leenkiat Soh](#) and [Adam Eck](#) (University of Nebraska-Lincoln)

CISER

(CED²AR): Comprehensive Extensible Data Documentation and Access Repository

- Method for solving the data curation problem that confronts the custodians of restricted-access research data and the scientific users of such data
- Accommodates physical security and access limitation protocols, and allows for much improved provenance tracking
- Metadata repository system that allows researchers to search, browse, access, and cite confidential data and metadata (via a web-based user interface or programmatically through a search API)

Proposed a <dataAccs> Solution at EDDI12 in Bergen

NCRN DDI Solution at the Variable Level: <dataAccs>

```
<studyDescr>
  <citation> [8 lines]
  <dataAccs ID="A1">
    <useStmt>
      <conditions>Public</conditions>
    </useStmt>
  </dataAccs>
  <dataAccs ID="A2">
    <useStmt>
      <confDec>To download this dataset, the user must obtain Special Sworn Status from the United States Census Bureau.</confDec>
      <conditions>Confidential</conditions>
    </useStmt>
  </dataAccs>
  <dataAccs ID="A3">
    <useStmt>
      <confDec>You're never gonna see this data.</confDec>
      <conditions>Need to know</conditions>
    </useStmt>
  </dataAccs>
</studyDescr>
```

Variable Level Solution (continued)

```
<var ID="V1500" dcml="0" files="F3" intrvl="discrete" name="totfam_kids" access="A1">
  <location width="12"/>
  <labl>Total Number of Children in Family</labl>
  <valrng> [2 lines]
  <sumStat type="vald">1000</sumStat>
  <sumStat type="invd">0</sumStat>
  <catgry> [3 lines]
  <catgry> [3 lines]
  <catgry> [3 lines]
  <catgry> [3 lines]
  <catgry> [3 lines]
  <catgry> [3 lines]
  <catgry> [3 lines]
  <catgry> [3 lines]
  <varFormat schema="other" type="numeric"/>
</var>
<var ID="V1588" dcml="0" files="F3" intrvl="contin" name="totinc" access="A2">
  <location width="12"/>
  <labl>Total Personal Income</labl>
  <valrng> [2 lines]
  <sumStat type="vald">240</sumStat>
  <sumStat type="invd">760</sumStat>
  <sumStat type="min">-278.739</sumStat>
  <sumStat type="max">39515.631</sumStat>
  <sumStat type="mean">1861.779</sumStat>
  <sumStat type="stdev">4015.033</sumStat>
  <varFormat schema="other" type="numeric"/>
</var>
```

No DDI Solution at the level of a *Value Label*

```
<var ID="V1588" dcml="0" files="F3" intrvl="contin" name="totinc" access="A1">
  <location width="12"/>
  <labl>Total Personal Income</labl>
  <catgry>
    <catValu>0</catValu>
    <labl>5-25k</labl>
  </catgry>
  <catgry>
    <catValu>1</catValu>
    <labl>25-75k</labl>
  </catgry>
  <catgry>
    <catValu>2</catValu>
    <labl>75-125k</labl>
  </catgry>
  <catgry>
    <catValu>3</catValu>
    <labl>125-250k</labl>
  </catgry>
  <catgry access="A2">
    <catValu>4</catValu>
    <labl>250k+</labl>
  </catgry>
  <varFormat schema="other" type="numeric"/>
</var>
```

Small tweak to the DDI Codebook Schema would fix this.

<dataAccs> developments since EDDI12

- In Lagoze, Block et.al. (2013) we more completely described the solution for embedding field-specific and value-specific cloaking in DDI Metadata*
- Proposed formal change to DDI 2.5 (April 2013)
- Brought modified “DDI 2.5.NCRN” schema online for testing (Fall 2013)
- Look forward to DDI Technical Implementation Committee taking up our proposal

*Lagoze, C., Block, W., Williams, J., Abowd, J. M., & Vilhuber, L. (2013). Data Management of Confidential Data. In *International Data Curation Conference*. Amsterdam.

Select Cornell NCRN Publications

Forthcoming. “Lagoze, Carl, Lars Vilhuber, Jeremy Williams, Benjamin Perry, and William C. Block, “CED²AR: The Comprehensive Extensible Data Documentation and Access Repository.” In: Proceedings of the ACM/IEEE Joint Conference on Digital Libraries (JCDL), London UK, September 2014.

2013 Lagoze, Carl, with William C. Block, Jeremy Williams, John M. Abowd, and Lars Vilhuber. “Data Management of Confidential Data”. In: International Journal of Digital Curation 8.1, pp.265-278. DOI: 10.2218/ijdc.v8il.259

2012 Abowd, John M., Lars Vilhuber, and William C. Block. “A Proposed Solution to the Archiving and Curation of Confidential Scientific Inputs.” In: Privacy in Statistical Databases. Ed. By Josep Domingo-Ferrer and Ilenia Tinnirello. Vol. 7556. Lecture Notes in Computer Science. Springer, pp.216-225. DOI: 10.1007/978-3-642-33627-0_17

Provenance

“data provenance, one kind of metadata, pertains to the derivation history of a data product starting from its original sources” [...] “from it, one can ascertain the quality of the data base and its ancestral data and derivations, track back sources of errors, allow automated reenactment of derivations to update the data, and provide attribution of data sources”*

*Simmhan, Plale, and Gannon, “A survey of data provenance in e-science,” ACM Sigmod Record, 2005

Provenance and Metadata

Not (currently) a “native” component of DDI, closest thing is:

```
<xs:complexType name="othrStdyMatType">
  <xs:complexContent>
    <xs:extension base="baseElementType">
      <xs:sequence>
        <xs:element ref="relMat" minOccurs="0" maxOccurs="unbounded"/>
        <xs:element ref="relStdy" minOccurs="0" maxOccurs="unbounded"/>
        <xs:element ref="relPubl" minOccurs="0" maxOccurs="unbounded"/>
        <xs:element ref="othRefs" minOccurs="0" maxOccurs="unbounded"/>
      </xs:sequence>
    </xs:extension>
  </xs:complexContent>
</xs:complexType>
```

Downside: No structure. Mostly verbose entries.

2013 work with PROV

- Explored encoding PROV in RDF/XML* (Required use of CDATA tag to avoid interfering with schema compliance; deemed less promising)
- More recently: exploring W3C PROV Model as basis for encoding provenance metadata in DDI

W3C PROV Model is based upon:

- **entities** that are physical, digital, and conceptual things in the world;
- **activities** that are dynamic aspects of the world that change and create entities; and
- **agents** that are responsible for activities.
- A set of **relationships** that can exist between them that express attribution, delegation, derivation, etc.

*Lagoze, C., Williams, J., & Vilhuber, L. (2013). Encoding Provenance Metadata for Social Science Datasets. In *7th Metadata and Semantics Research Conference*. Thessaloniki.

CISER

The American Community Survey (ACS)

- Ongoing statistical survey conducted by the U.S. Census Bureau
- Approximately 250,000 surveys/month (3 million per year)
- Replacement for detailed long-form decennial census

The image shows a screenshot of a web browser displaying the American Community Survey questionnaire booklet. The browser address bar shows the URL: www.census.gov/acs/www/Downloads/questionnaires/2014/Quest14.pdf. The booklet header includes the U.S. Department of Commerce logo and the text: "U.S. DEPARTMENT OF COMMERCE Economics and Statistics Administration U.S. CENSUS BUREAU". The title of the survey is "THE American Community Survey". A green box highlights the text: "This booklet shows the content of the American Community Survey questionnaire." The "Start Here" section provides instructions on how to respond, either online at <https://respond.census.gov/acs> or by mail. It also includes a section for help, with the phone number 1-800-354-7271. The "How many people are living or staying at this address?" section includes instructions to include everyone living or staying here for more than 2 months, including oneself and anyone else staying here who does not have another place to stay, even if they are here for 2 months or less. It also states to do not include anyone who is living somewhere else for more than 2 months, such as a college student living away or someone in the Armed Forces on deployment. The "Number of people" field is currently empty.

13194014

U.S. DEPARTMENT OF COMMERCE
Economics and Statistics Administration
U.S. CENSUS BUREAU

THE American Community Survey

This booklet shows the content of the American Community Survey questionnaire.

Start Here

Respond online today at:
<https://respond.census.gov/acs>
OR
Complete this form and mail it back as soon as possible.

This form asks for information about the people who are living or staying at the address on the mailing label and about the house, apartment, or mobile home located at the address on the mailing label.

If you need help or have questions about completing this form, please call 1-800-354-7271. The telephone call is free.

Telephone Device for the Deaf (TDD):
Call 1-800-582-8330. The telephone call is free.

¿NECESITA AYUDA? Si usted habla español y necesita ayuda para completar su cuestionario, llame sin cargo alguno al 1-877-833-5625. Usted también puede completar su entrevista por teléfono con un entrevistador que habla español. O puede responder por Internet en: <https://respond.census.gov/acs>

For more information about the American

Please print today's date.
Month Day Year

Please print the name and telephone number of the person who is filling out this form. We may contact you if there is a question.

Last Name

First Name MI

Area Code + Number

How many people are living or staying at this address?
• INCLUDE everyone who is living or staying here for more than 2 months.
• INCLUDE yourself if you are living here for more than 2 months.
• INCLUDE anyone else staying here who does not have another place to stay, even if they are here for 2 months or less.
• DO NOT INCLUDE anyone who is living somewhere else for more than 2 months, such as a college student living away or someone in the Armed Forces on deployment.

Number of people

ACS Question on Ancestry or Ethnic Origin

13194089

Person 1

→ Please copy the name of Person 1 from page 2, then continue answering questions below.

Last Name

First Name MI

7 Where was this person born?

In the United States – Print name of state.

Outside the United States – Print name of foreign country, or Puerto Rico, Guam, etc.

8 Is this person a citizen of the United States?

Yes, born in the United States → SKIP to question 10a

Yes, born in Puerto Rico, Guam, the U.S. Virgin Islands, or Northern Marianas

Yes, born abroad of U.S. citizen parent or parents

Yes, U.S. citizen by naturalization – Print year of naturalization

No, not a U.S. citizen

9 When did this person come to live in the United States? Print numbers in boxes.

Year

10 a. At any time IN THE LAST 3 MONTHS, has this person attended school or college? Include only nursery or preschool, kindergarten, elementary school, home school, and schooling which leads to a high school diploma or a college degree.

No, has not attended in the last 3 months → SKIP to question 11

Yes, public school, public college

Yes, private school, private college

11 What is the highest degree or level of school this person has COMPLETED? Mark (X) ONE box. If currently enrolled, mark the previous grade or highest degree received.

NO SCHOOLING COMPLETED

No schooling completed

NURSERY OR PRESCHOOL THROUGH GRADE 12

Nursery school

Kindergarten

Grade 1 through 11 – Specify grade 1 – 11

12th grade – **NO DIPLOMA**

HIGH SCHOOL GRADUATE

Regular high school diploma

GED or alternative credential

COLLEGE OR SOME COLLEGE

Some college credit, but less than 1 year of college credit

1 or more years of college credit, no degree

Associate's degree (for example: AA, AS)

Bachelor's degree (for example: BA, BS)

AFTER BACHELOR'S DEGREE

Master's degree (for example: MA, MS, MEng, MEd, MSW, MBA)

Professional degree beyond a bachelor's degree (for example: MD, DDS, DVM, LLB, JD)

Doctorate degree (for example: PhD, EdD)

F Answer question 12 if this person has a bachelor's degree or higher. Otherwise, SKIP to question 13.

12 This question focuses on this person's BACHELOR'S DEGREE. Please print below the specific major(s) of any BACHELOR'S DEGREES this person has received. (For example: chemical engineering, elementary teacher education, organizational psychology)

13 What is this person's ancestry or ethnic origin?

(For example: Italian, Jamaican, African Am., Cambodian, Cape Verdean, Norwegian, Dominican, French Canadian, Haitian, Korean, Lebanese, Polish, Nigerian, Mexican, Taiwanese, Ukrainian, and so on.)

14 a. Does this person speak a language other than English at home?

Yes

No → SKIP to question 15a

b. What is this language?

(For example: Korean, Italian, Spanish, Vietnamese)

c. How well does this person speak English?

Very well

Well

Not well

Not at all

15 a. Did this person live in this house or apartment 1 year ago?

Person is under 1 year old → SKIP to question 16

Yes, this house → SKIP to question 16

No, outside the United States and Puerto Rico – Print name of foreign country, or U.S. Virgin Islands, Guam, etc., below; then SKIP to question 16

No, different house in the United States or Puerto Rico

b. Where did this person live 1 year ago?

Address (Number and street name)

CISER

Three Use Cases: Researchers interested in people of Alsatian, Andorran, and Cypriot Ancestry

<u>Code</u>	<u>Write-In</u>
001-099	WESTERN EUROPE (EXCEPT SPAIN)
001	ALSATIAN
002	ANDORRAN
003	AUSTRIAN
004	TIROL
005	BASQUE
006	FRENCH BASQUE
007	SPANISH BASQUE
008	BELGIAN
009	FLEMISH
010	WALLOON
011	BRITISH
012	BRITISH ISLES
013	CHANNEL ISLANDER
014	GIBRALTAR
015	CORNISH
016	CORSICAN
017	CYPRIOT
018	GREEK CYPRIOTE
019	TURKISH CYPRIOTE
020	DANISH

- U.S. Census Bureau Documentation
- Ancestry Code List
- 2012 ACS

CISER

Multiple Sources of Data originating from the ACS: Examples of Aggregate Data

U.S. Department of Commerce
United States Census Bureau
AMERICAN FactFinder
MAIN | COMMUNITY FACTS | GUIDED SEARCH | ADVANCED SEARCH | DOWNLOAD CENTER

NHGIS data finder
Filter » Options » Review

Apply Filters [How to use the data finder \(pdf\)](#)

Geographic Levels nation
Years 2012

12 fx Estimates: Alsatian

	A	B	C	D	E	F	G	H	I
1	GISJOIN	YEAR	NATION	NATIONA	NAME	OJ2E001	OJ2E002	OJ2E003	OJ2E004
2	GIS Join Match Code	Data File Year	Nation Name	Nation Code	Area Name	Estimates: Total	Estimates: Afghan	Estimates: Albanian	Estimates: Alsatian
3	G1	2012	United State	1	United State	313914040	94244	185696	6626

2012 ACS 1-year Estimate: 6,626 individuals of Alsatian Ancestry living in the United States

CISER

Multiple Sources of Data originating from the ACS: Example of PUMS Microdata

The screenshot shows an Excel spreadsheet titled '2012 PUMS Ancestry Code List'. The table has four columns: 'PUMS Code', 'Ancestry Description', 'Ancestry Code', and 'Corresponding Detailed Ancestry Code'. The data lists various ancestry categories and their corresponding codes. A red arrow points to the 'Ancestry.xlsx' file in the background, and another red arrow points to the 'Ancestry Code' column in the table.

	A	B	C	D
		2012 PUMS Ancestry Code List		
	PUMS Code	Ancestry Description	Ancestry Code	Corresponding Detailed Ancestry Code
2		001	Alsatian	001 Alsatian
3		003	Austrian	003 Austrian
4				004 Tirol
5		005	Basque	005 Basque
6				006 French Basque
7				007 Spanish Basque
8		008	Belgian	008 Belgian
9				010 Walloon
10		009	Flemish	009 Flemish
11		011	British	011 British
12				013 Channel Islander
13				014 Gibraltar
14		012	British Isles	012 British Isles
15		020	Danish	020 Danish
16				023 Faroe Islander
17		021	Dutch	021 Dutch
18				029 Frisian
19		022	English	015 Cornish
20				022 English
21		024	Finnish	024 Finnish
22				025 Karelian
23		026	French	016 Corsican
24				026 French
25				027 Lorraine
26				028 Breton
27				083 Occitan
28		032	German	032 German
29				033 Bavaria
30				

ACS 2012 PUMS:
ANCESTRY Code is
001 for Alsatian

Multiple Sources of Data originating from the ACS: Example of IPUMS-USA

IPUMS-USA: descr: ANCESTR1

https://usa.ipums.org/usa-action/variables/ANCESTR1#codes_section

MINNESOTA POPULATION CENTER, UNIVERSITY OF MINNESOTA

IPUMS USA

Home Select Data FAQ Help Login

ANCESTR1 Remove from cart Change samples

Ancestry, first response
Group: [Race, Ethnicity, and Nativity — PERSON](#)

Description Comparability Universe Codes Availability Questionnaire Text Flags

Codes and Frequencies

Category availability view General codes
 Case-count view Detailed codes

Code	Label	2012
	WESTERN EUROPE (EXCEPT SPAIN)	.
001	Alsatian, Alsace-Lorraine	75
002	Andorran	.
003	Austrian	4,172
004	Tirolean	.
	Basque:	.
005	Basque	482
006	French Basque	.
008	Belgian	2,444
009	Flemish	66
010	Walloon	.
011	British	10,077
012	British Isles	537
013	Channel Islander	.
014	Gibraltar	.
015	Cornish	.
016	Corsican	.
017	Cypriot	.
018	Greek Cypriote	.
019	Turkish Cypriote	.
020	Danish	8,412
021	Dutch	26,301
022	English	191,812
023	Faeroe Islander	.
024	Finnish	4,963

IPUMS-USA for ACS 2012:

- 001 Alsatian ANCESTRY Code
- 75 cases in the sample

CISER

- Let's review...

	2012 ACS Code List	ACS 2012 PUMS	IPUMS-USA	AFF	NHGIS	
Alsatian	YES (001)	YES (001)	YES (75 cases)	6,626 (est.)	6,626 (est.)	
Andorran						
Cypriots						

CISER

Three Use Cases: Researchers interested in people of Alsatian, Andorran, and Cypriot Ancestry

Ancestry Code List

<u>Code</u>	<u>Write-In</u>
001-099	WESTERN EUROPE (EXCEPT SPAIN)
001	ALSATIAN
002	ANDORRAN
003	AUSTRIAN
004	TIROL
005	BASQUE
006	FRENCH BASQUE
007	SPANISH BASQUE
008	BELGIAN
009	FLEMISH
010	WALLOON
011	BRITISH
012	BRITISH ISLES
013	CHANNEL ISLANDER
014	GIBRALTAR
015	CORNISH
016	CORSICAN
017	CYPRIOT
018	GREEK CYPRIOTE
019	TURKISH CYPRIOTE
020	DANISH

- U.S. Census Bureau Documentation
- Ancestry Code List
- 2012 ACS

CISER

	2012 ACS Code List	ACS 2012 PUMS	IPUMS-USA	AFF	NHGIS	
Alsatian	YES (001)	YES (001)	YES (75 cases)	6,626 (est.)	6,626 (est.)	
Andorran	YES (002)					
Cypriots	YES (017)					

CISER

U.S. Department of Commerce

United States[™]
Census
Bureau

AMERICAN
FactFinder 

Basque	46,874	+/-4,464
Belgian	230,104	+/-8,312
Brazilian	321,544	+/-14,571
British	931,514	+/-17,814

NHGIS data finder

AD	AE	AF	AG	AH	AI
OJ2E025	OJ2E026	OJ2E027	OJ2E028	OJ2E029	OJ2E030
Estimates: Cajun	Estimates: Canadian	Estimates: Carpatho Rusyn	Estimates: Celtic	Estimates: Croatian	Estimates: Cypriot
96349	509952	4944	34844	276808	6486

CISER

	2012 ACS Code List	ACS 2012 PUMS	IPUMS-USA	AFF	NHGIS	
Alsatian	YES (001)	YES (001)	YES (75 cases)	6,626 (est.)	6,626 (est.)	
Andorran	YES (002)					
Cypriots	YES (017)			6,486 (est.)	6,486 (est.)	

CISER

MINNESOTA POPULATION CENTER, UNIVERSITY OF MINNESOTA

IPUMS USA

1			
2	Code	Label	2012
3			acs
4		WESTERN EUROPE (EXCEPT SPAIN)	.
5	001	Alsatian, Alsace-Lorraine	X
6	002	Andorran	.
7	003	Austrian	X
8	004	Tirolean	.
9		Basque:	.
10	005	Basque	X
11	006	French Basque	.
12	008	Belgian	X
13	009	Flemish	X
14	010	Walloon	.
15	011	British	X
16	012	British Isles	X
17	013	Channel Islander	.
18	014	Gibraltar	.
19	015	Cornish	.
20	016	Corsican	.
21	017	Cypriot	.
22	018	Greek Cypriote	.
	019	Turkish Cypriote	.

CISER

	2012 ACS Code List	ACS 2012 PUMS	IPUMS- USA	AFF	NHGIS	
Alsatian	YES (001)	YES (001)	YES (75 cases)	6,626 (est.)	6,626 (est.)	
Andorran	YES (002)					
Cypriots	YES (017)	NO	NO	6,486 (est.)	6,486 (est.)	

CISER

Three Use Cases: Researchers interested in people of Alsatian, Andorran, and Cypriot Ancestry

Ancestry Code List

<u>Code</u>	<u>Write-In</u>
001-099	WESTERN EUROPE (EXCEPT SPAIN)
001	ALSATIAN
002	ANDORRAN
003	AUSTRIAN
004	TIROL
005	BASQUE
006	FRENCH BASQUE
007	SPANISH BASQUE
008	BELGIAN
009	FLEMISH
010	WALLOON
011	BRITISH
012	BRITISH ISLES
013	CHANNEL ISLANDER
014	GIBRALTAR
015	CORNISH
016	CORSICAN
017	CYPRIOT
018	GREEK CYPRIOTE
019	TURKISH CYPRIOTE
020	DANISH

- U.S. Census Bureau Documentation
- Ancestry Code List
- 2012 ACS

CISER

IPUMS USA

AMERICAN FactFinder



ANCESTR1

Ancestry, first resp

Group: [Race, Etc](#)

Description

1

United States

Estimate Margin of Error

09	Total:	313,914,040	*****
09	Afghan	94,244	+/-8,662
	Albanian	185,696	+/-13,225
	Alsatian	6,626	+/-1,186
	American	21,026,438	+/-90,090
	Arab:	1,524,666	+/-34,383
	Egyptian	214,890	+/-13,737

An 'X' indicate

Code Label

2012
acs

	WESTERN EUROPE (EXCEPT SPAIN)	.
001	Alsatian, Alsace-Lorraine	X
002	Andorran	.
003	Austrian	X

CISER

	2012 ACS Code List	ACS 2012 PUMS	IPUMS-USA	AFF	NHGIS	
Alsatian	YES (001)	YES (001)	YES (75 cases)	6,626 (est.)	6,626 (est.)	
Andorran	YES (002)	NO	NO	NO	NO	
Cypriots	YES (017)	NO	NO	6,486 (est.)	6,486 (est.)	

CISER

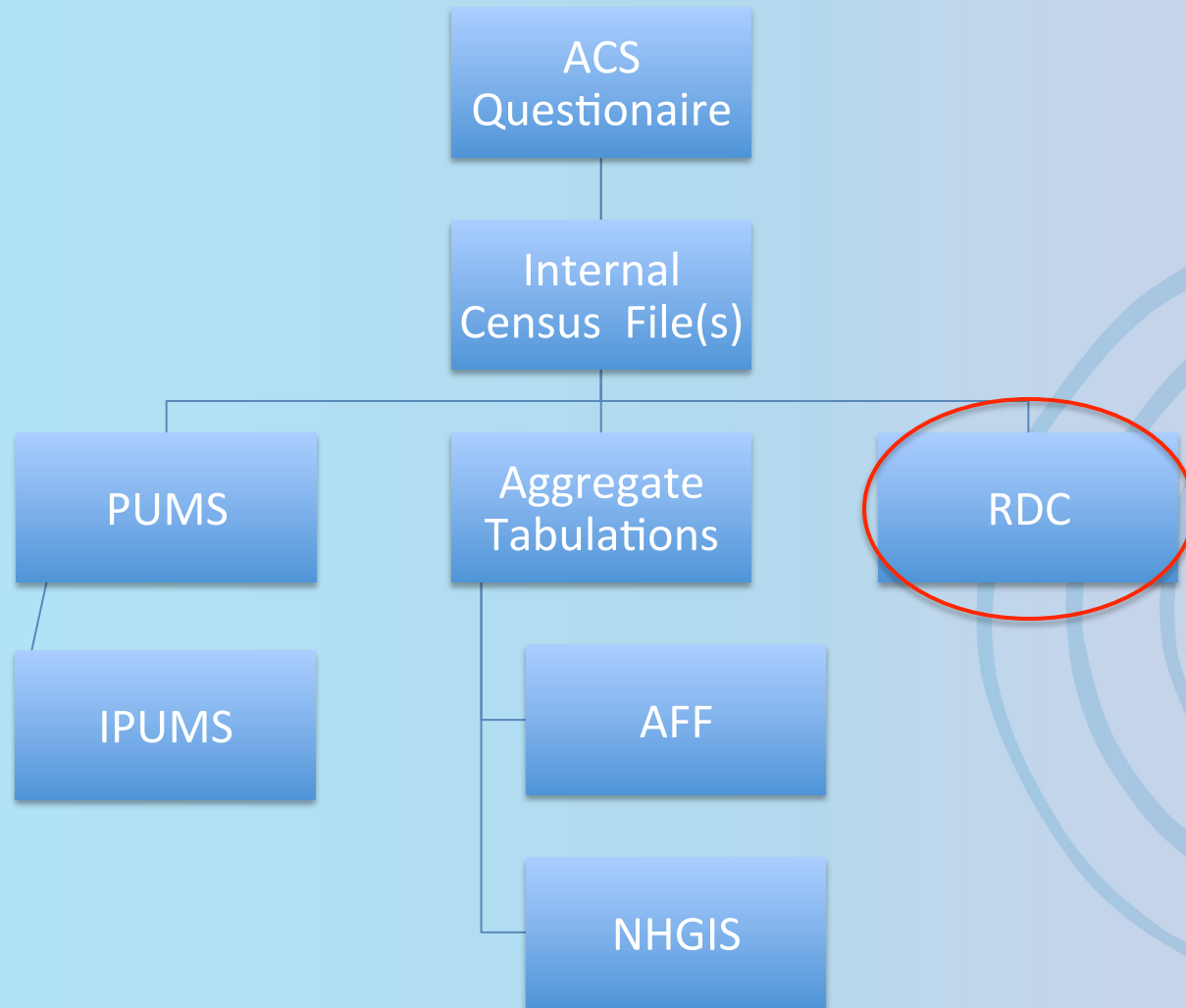
Three Use Cases: Researchers interested in people of Alsatian, Andorran, and Cypriot Ancestry

2012 PUMS Ancestry Code List			
...
995	Mixture	995	Mixture
996	Uncodable Entries	996	Uncodable Entries
997	Other Groups	002	Andorran
		017	Cypriot
		018	Greek Cypriote
		019	Turkish Cypriote
		020	Lapp
		021	Liechtensteiner
		022	Manx
		023	Monegasque
		101	Azerbaijani
		107	Ruthenian
		108	Cossack
		117	Finno Ugrian
		118	Mordovian
		119	Voytak
		120	Gruzia
		127	Kalmyk

- U.S. Census Bureau Documentation
- Ancestry Code List
- 2012 ACS

CISER

- Simple Provenance of ACS Data Files

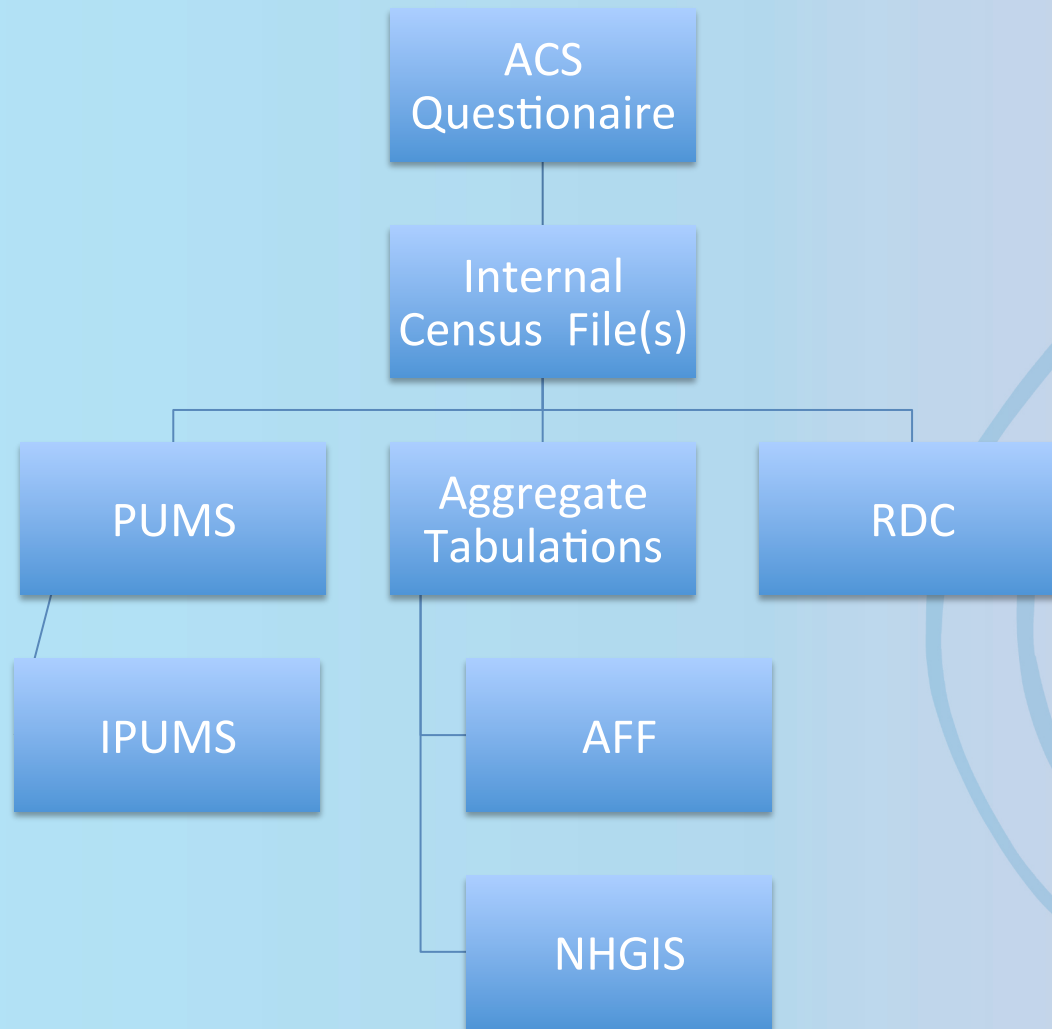


CISER

	2012 ACS Code List	ACS 2012 PUMS	IPUMS-USA	AFF	NHGIS	RDC
Alsatian	YES (001)	YES (001)	YES (75 cases)	6,626 (est.)	6,626 (est.)	Yes
Andorran	YES (002)	NO	NO	NO	NO	?
Cypriots	YES (017)	NO	NO	6,486 (est.)	6,486 (est.)	?

CISER

- Provenance of ACS Data Files



PROV at the Dataset/File Level

- Mentioned earlier our exploration of encoding PROV in RDF/XML* (CDATA; less promising)
- More recently: exploring W3C PROV Model as basis for encoding provenance metadata in DDI XML using <relStdy>
 - “...information on the relationship of the current data collection to others (e.g., predecessors, successors, other waves or rounds or to other editions of the same file). This would include the names of additional data collections generated from the same data collection vehicle plus other collections directed at the same general topic, which can take the form of bibliographic citations.”
- Others working to develop an RDF encoding for DDI metadata that could easily accommodate RDF-encoding of provenance metadata
- Both solutions might be viable; implementation could depend on local preferences

Embed PROV within <RelStdy>

Step 1 - Material Reference Complex Type with Prov

To include prov:Document within the <relStdy> element, a new complex type called 'materialReferenceWithProvType', which inherits from materialReferenceType can be introduced as follows:

```
<xs:complexType name="materialReferenceWithProvType" mixed="true">
  <xs:complexContent>
    <xs:extension base="materialReferenceType">
      <xs:sequence>
        <xs:element ref="prov:document" minOccurs="0" maxOccurs="1" />
      </xs:sequence>
      <xs:attribute name="provDocURI" type="xs:anyURI" use="optional"/>
    </xs:extension>
  </xs:complexContent>
</xs:complexType>
```

This allows PROV document to be embedded or referenced by URI.

Embed PROV within <RelStdy> (Cont)

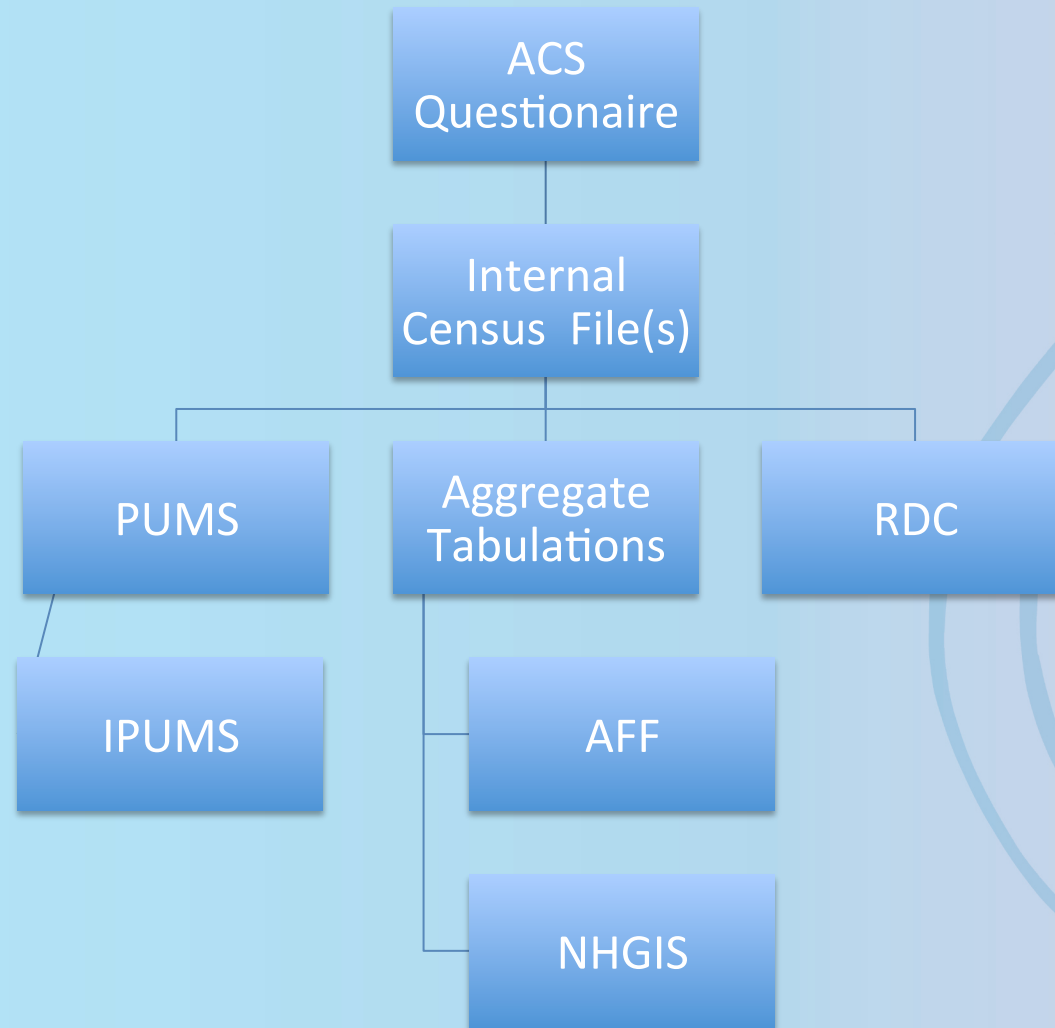
Step 2 - Modify the type of relStdy to the new complex type

The relStdy element is changed to inherit from this 'materialReferenceWithProvType', and examples are given

```
<xs:element name="relStdy" type="materialReferenceWithProvType">
  <xs:annotation>
    <xs:documentation>
      <xhtml:div>
        <xhtml:h1 class="element_title">Related Studies</xhtml:h1>
        <xhtml:div>
          <xhtml:h2 class="section_header">Description</xhtml:h2>
          <xhtml:div class="description">Information on the relationship of the current data collection to others [4 lines]
        </xhtml:div>
        <xhtml:div> [8 lines]
        <xhtml:div> [3 lines]
        <xhtml:div>
          <xhtml:h2 class="section_header">Prov Example - Embedded Provenance Document</xhtml:h2>
          <xhtml:div class="example">
            <xhtml:samp class="xml_sample"><![CDATA[
              <relStdy>ICPSR distributes a companion study to this collection titled FEMALE LABOR
              FORCE PARTICIPATION AND MARITAL INSTABILITY, 1980: [UNITED STATES] (ICPSR 9199).
              <prov:document>...</prov:document></relStdy>
            ]]></xhtml:samp>
          </xhtml:div>
        </xhtml:div>
        <xhtml:div>
          <xhtml:h2 class="section_header">Prov Example - Referenced Provenance Document</xhtml:h2>
          <xhtml:div class="example">
            <xhtml:samp class="xml_sample"><![CDATA[
              <relStdy provDocURI="http://dx.doi.org/10.1234/exProv123">ICPSR distributes a companion
              study to this collection titled FEMALE LABOR FORCE PARTICIPATION AND MARITAL INSTABILITY,
              1980: [UNITED STATES] (ICPSR 9199).</relStdy>
            ]]></xhtml:samp>
          </xhtml:div>
        </xhtml:div>
      </xs:documentation>
    </xs:annotation>
  </xs:element>
```

CISER

- Provenance of ACS Data Files



Variable Level Provenance

A single attribute is added to the variable type

```
</xs:attribute/>
<xs:attribute name="geoVocab" type="xs:string"/>
<xs:attribute name="catQty" type="xs:string"/>
<xs:attribute name="representationType">
  <xs:simpleType>
    <xs:restriction base="xs:NMTOKEN">
      <xs:enumeration value="text"/>
      <xs:enumeration value="numeric"/>
      <xs:enumeration value="code"/>
      <xs:enumeration value="datetime"/>
      <xs:enumeration value="other"/>
    </xs:restriction>
  </xs:simpleType>
  <xs:attribute ref="prov:ref" use="optional"/>
  <xs:attribute name="otherRepresentationType" type="xs:NMTOKEN" use="optional"/>
</xs:extension>
</xs:complexContent>
</xs:complexType>
```

Variable Level Provenance (cont.)

In our implementation, a variable would reference a prov:Bundle that would be found within embedded prov:Document. Here is an example of a prov:Bundle:

```
<prov:bundle prov:id="ex:bundle1">
  <ex:version>1</ex:version>
</prov:bundle>

<prov:bundleContent prov:id="ex:bundle1">
  <prov:entity prov:id="ex:report1"/>

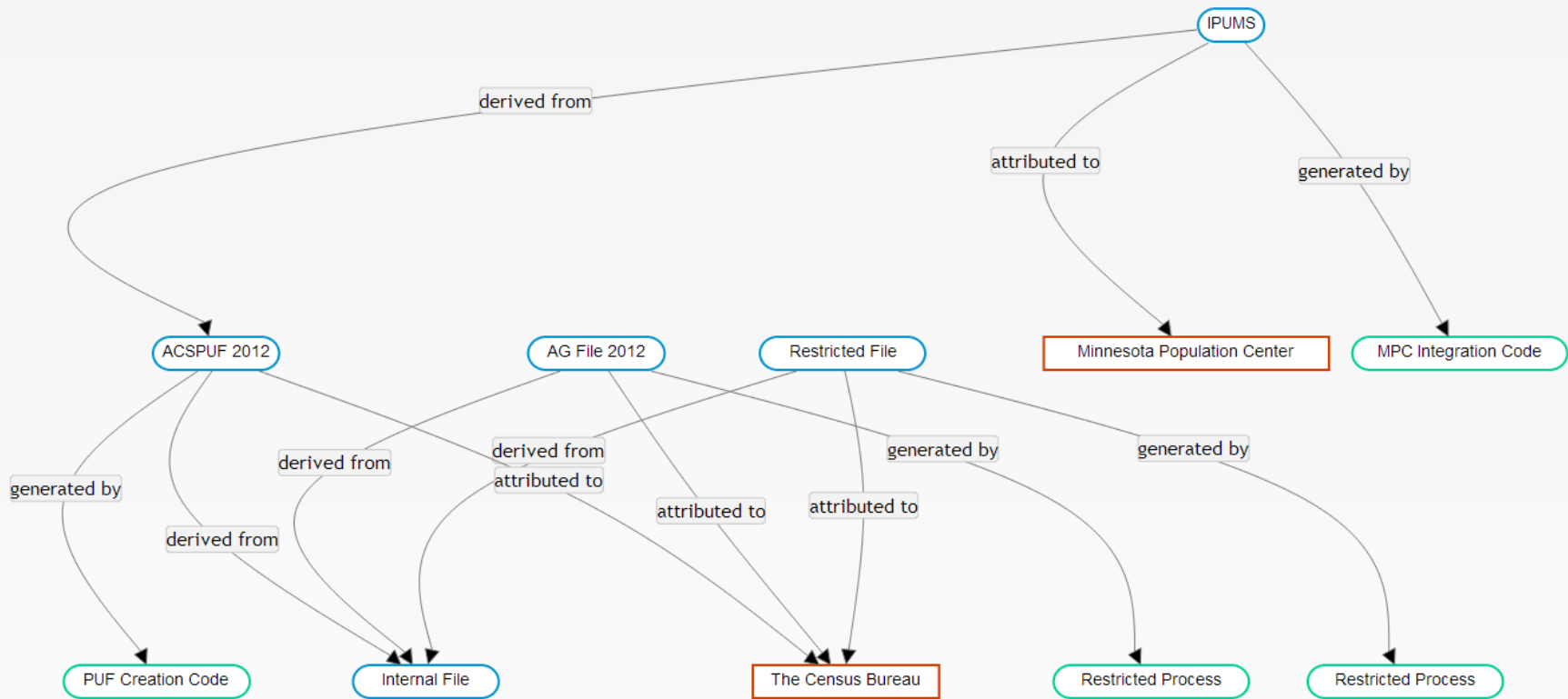
  <prov:entity prov:id="ex:report2">
    <prov:type xsi:type="xsd:QName">report</prov:type>
    <ex:version>2</ex:version>
  </prov:entity>

  <prov:wasGeneratedBy>
    <prov:entity prov:ref="ex:report2"/>
    <prov:time>2012-05-25T11:00:01</prov:time>
  </prov:wasGeneratedBy>

  <prov:wasDerivedFrom>
    <prov:generatedEntity prov:ref="ex:report2"/>
    <prov:usedEntity prov:ref="ex:report1"/>
  </prov:wasDerivedFrom>
</prov:bundleContent>
```

CISER

Prov



NCRN Meeting Spring 2014

[View](#)[Register](#)

The eight nodes of the NSF-Census Research Network (NCRN) will hold a technical mini-symposium on Thursday and Friday, May 22-23, 2014. The program will comprise two one-half day sessions focused on topics of broad interest to the nodes, the Census Bureau and the federal statistical system. Presentations will highlight research performed by the NCRN nodes, with opportunity for questions and comments from agency researchers.

Thursday, May 22

9:00-9:30 Opening remarks by the Director of the Census Bureau, John Thompson

9:30-12:00 **Session 1: Data Documentation Initiative (DDI) Metadata within the Federal Statistical System: Implementation Challenges, and Provenance Encoding**

DDI Presentations at NCRN Spring Meeting

Data Documentation Initiative (DDI) Metadata within the Federal Statistical System: Implementation Challenges, Record Linkage, and Provenance Encoding

Jay Greenfield and Sophia Kuan (Booz Allen Hamilton): Describing Adaptive/Responsive Protocols and Data Processing in DDI

Tim Mulcahy (NORC): DDI and Record Linkage

Jeremy Williams (Cornell University): Encoding Provenance Metadata for Social Science Datasets

DDI Tools Demonstration

Abdul K. Rahim and Pascal Heus (Metadata Technologies North America)

CISER

Cornell Institute for Social and Economic Research
A LEADER IN SOCIAL SCIENCE DATA AND COMPUTING



Thank you!

Questions?



block@cornell.edu

ncrn.cornell.edu



Cornell University