**CIVIC**
**DIGITAL**
**FELLOWSHIP**

# Novel Data Linkage in Support of the Decennial Census Digitization and Linkage Project
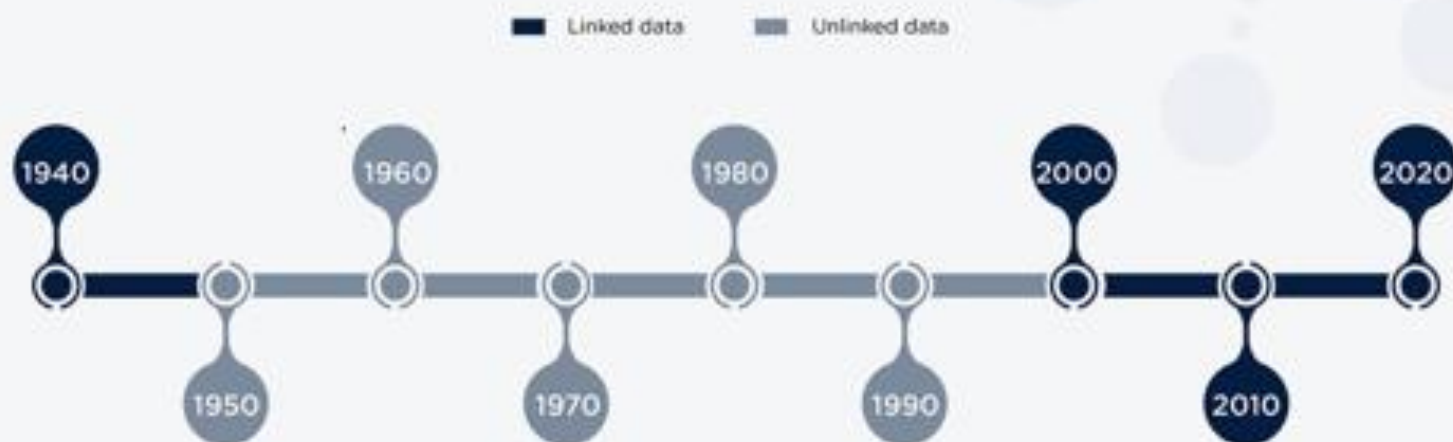
Sohail Kamdar

Supervised by Katie Genadek and John Sullivan

Economic Reimbursable Surveys Division

# Background and Motivations



**Linked Decennial Census Data**

Legend: ■ Linked data  ■ Unlinked data

1940 · 1950 · 1960 · 1970 · 1980 · 1990 · 2000 · 2010 · 2020

United States® Census Bureau

Linking person-level census microdata over time will fill a 5-decade gap in longitudinal infrastructure for researchers studying long term social, economic, and health outcomes.
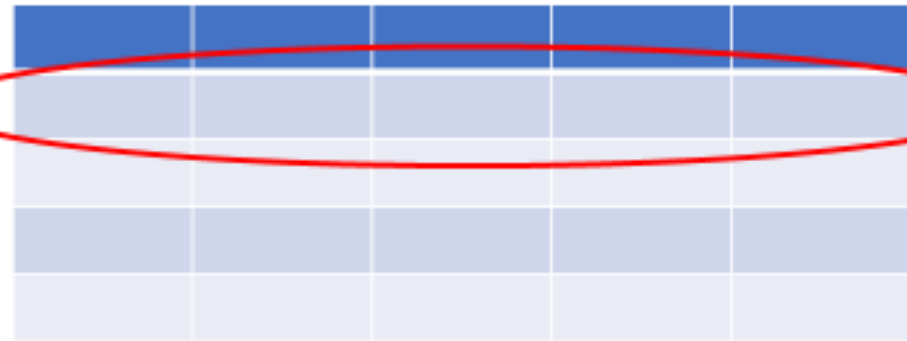
Names: Not currently captured

Bubbles: Currently in census microdata

# Objectives

## Linking keyed microfilm data to short form census microdata

1. Survey forms on microfilm converted to keyed datasets



2. Keyed image data and census micro data are linked

**Name Data**

**+**

Short form 100% file

# Linkage process:

| Name | Age | Number in Household | Sex | Relationship | County |
|---|---|---|---|---|---|
| Darth Vader | 45 | 3 | 0 | 0 | 0001 |
| Luke Skywalker | 25 | 3 | 0 | 2 | 0001 |
| Leia Organa | 26 | 3 | 1 | 2 | 0001 |

(Dataset from images)

1. Clean microdata and name data such that records are comparable

2. Match household data within geographic areas

3. Match person level data using demographic information

4. Use fuzzy matching to deal with data errors.

| Age | Number in Household | Sex | Relationship | County |
|---|---|---|---|---|
| 45 | 3 | 0 | 0 | 0001 |
| 25 | 3 | 0 | 2 | 0001 |
| 25 | 3 | 1 | 2 | 0001 |

(Census dataset)

United States® Census Bureau

5

# Additional Challenges and Potential solutions

**Challenges**

- Lack of string identifiers

- Opaque legacy documentation

- Geographic inconsistencies

**Solutions**

→ • Block data into useful variables

→ • Use only commonly coded variables

→ • Use alternative markers of geography to merge datasets

# Preliminary Results and Next Steps

- Match rates for keyed data and microdata are between 80-90% for sample geographies

- Errors and inconsistencies between datasets are documented

**Next steps:**

- Resolving challenges with missing data in rural areas

- Creating scalable solution to deal with differences within regions

United States® **Census** Bureau