# Leveraging *Local* Data Sources

May 20, 2016

Sallie Keller, Director

Virginia Tech
Biocomplexity Institute

SDAL
SOCIAL &
DECISION ANALYTICS
LABORATORY

# Data characterizes the world and the ways in which communication occurs between its individuals

## Infrastructure



- Condition
- Operations
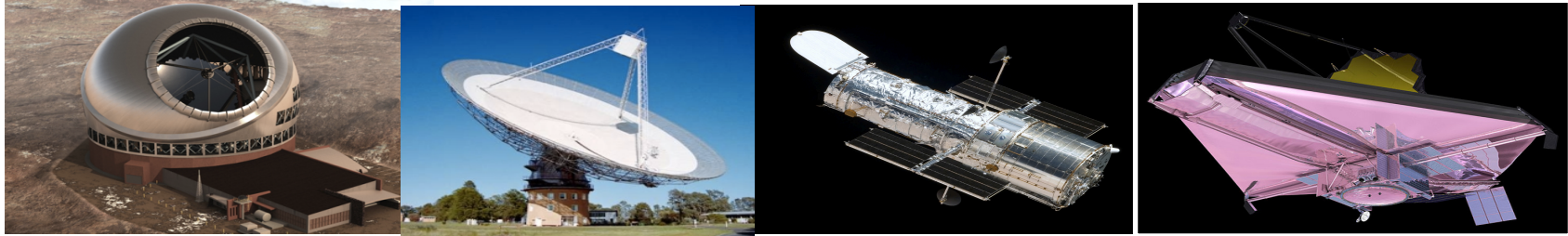- Resilience
- Sustainability

## Environment



- Climate
- Pollution
- Noise
- Flora/ Fauna

## People



- Relationships
- Location
- Economic Condition
- Communication
- Activities
- Health

# The New Lens for Social Behavioral Observing



- Data collected faster, while individuals are in the act of behaving in real life situation

- Adapt methods to make the best use of these data

- New data streams produce new discoveries but should not be allowed to degrade the scientific approach

# We are in the middle of an *ALL data* revolution

## Designed Data Collections

Statistically-designed and intentional observational data collections



## Administrative Data

Data collected for the administration of an organization or program



## Opportunity Data

Data generated as we move through our daily paces



## Procedural Data

Data derived from policies and procedures

# What about data quality?

## Traditional Approach:

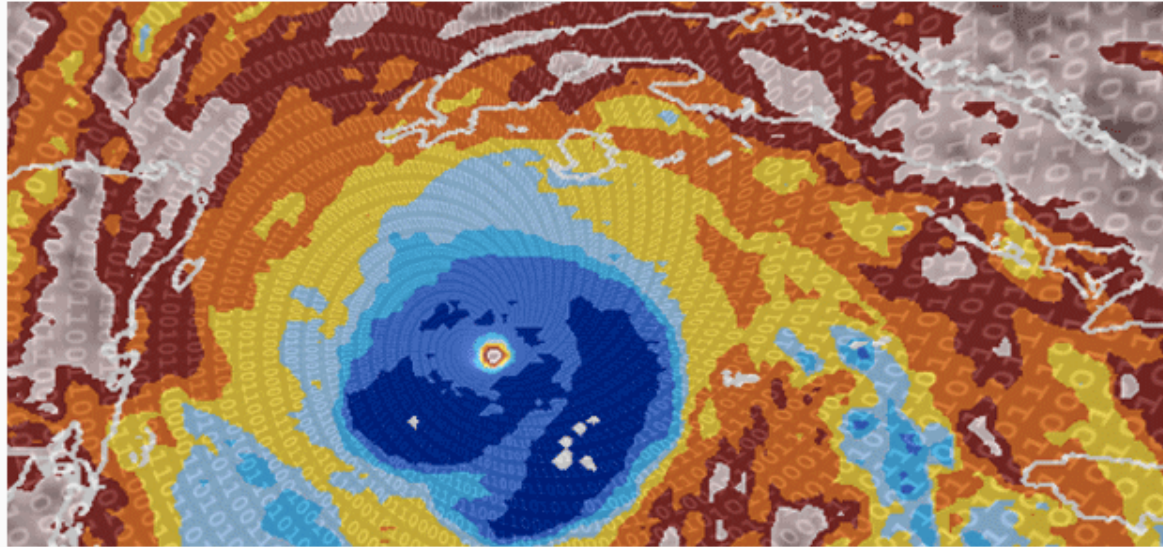- Control over measurement processes

- Control over collection processes
  - Optimization
  - information maximization

- Clear and controlled ownership

# Evolution of data quality by discipline

| Discipline | Contributions to Data Quality |
|---|---|
| Physical and Biological Sciences | • Experimental methods<br>• Data repositories/portals<br>• Reproducibility and replication |
| Engineering, IT, Business | • Pareto Principle (80 - 20 rule)<br>• Fitness-for-use<br>• Total Data Quality Management (TDQM)<br>• Data management<br>• Standards |
| Social and Behavioral Sciences | • Total Survey Error<br>• Randomized Control Trials, Observational Studies, and Natural Experiments |
| Statistics | • Decision theoretic approach<br>• Statistical methods<br>• Adoption of definitions, methods, and approaches by official statistics |

# In this data revolution, we have lost control!



5 keys to getting big data under control

By Paul McCloskey    May 20, 2013

# And then there is privacy!

Does Big Data Change the Privacy Landscape: A Review of the Issues

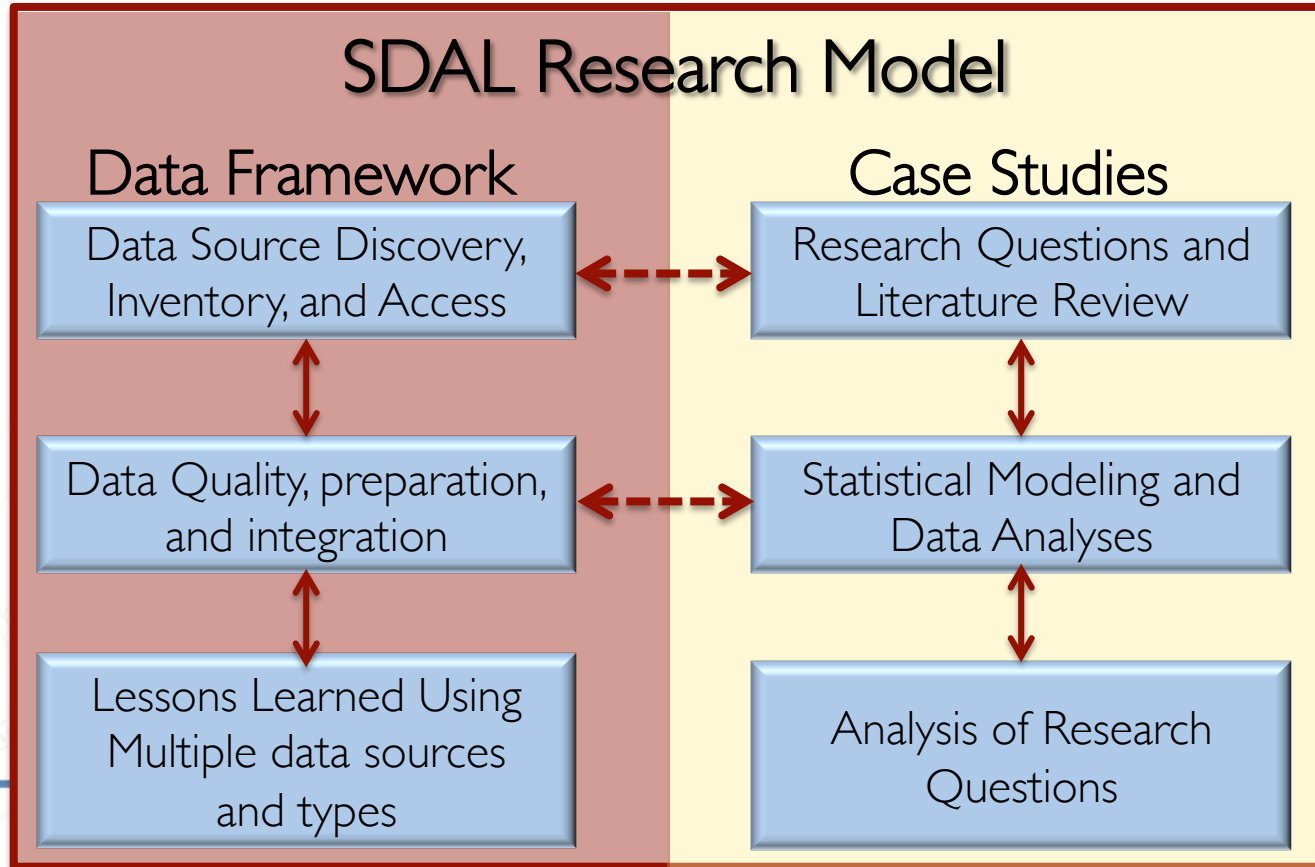Sallie Ann Keller, Stephanie Shipp, and Aaron Schroeder

Virginia Tech
Biocomplexity Institute

# Innovation data and measurement problems

- For what purpose?
  - Research, researcher access, policy, program evaluation, …
- What dimension or aspect of innovaion?
  - entrepreneurship, R&D, advances, productivity, productivity growth, technology, workforce, household innovators, automation, robotics, non-market diffusion, collective knowledge, …
- Can we find local data surrogates?

# A disciplined approach is needed to develop the theory and methods for using All the data

## SDAL Research Model

### Data Framework

| Data Source Discovery, Inventory, and Access |
| --- |

⇅

| Data Quality, preparation, and integration |
| --- |

⇅

| Lessons Learned Using Multiple data sources and types |
| --- |

### Case Studies

| Research Questions and Literature Review |
| --- |

⇅

| Statistical Modeling and Data Analyses |
| --- |

⇅

| Analysis of Research Questions |
| --- |

irginiaTech
complexity Institute
mplexity Institute

# Data Framework

**PROBLEM IDENTIFICATION:** Relevant Issues and Working Hypotheses

⇕

| LOCAL DATA<br>FEDERAL and STATE DATA<br>DESIGNED DATA<br>OPPORTUNISTIC DATA FLOWS | DATA SOURCES:<br>Discovery<br>Inventory<br>Acquisition |
| --- | --- |

⇕

**DATA STORAGE**

| ⇕ | ⇕ | ⇕ | ⇕ |
| --- | --- | --- | --- |
| DATA PROFILING | DATA PREPARATION | DATA LINKAGE | DATA EXPLORATION |
| ⇕ | ⇕ | ⇕ | ⇕ |

**MODELING AND ANALYSES**

⇕

**FITNESS-FOR-USE ASSESSMENT**

# Data inventory process in-action

- Purpose of organization collecting the data
- Description of the data
  - content
  - Unit represented
  - Longitudinal or cross-sectional
  - Geographic coverage
  - Timing of collection and release

- Data collection method
- Metadata and provenance
- Other issues
  - Selectivity
  - Stability
  - Accessibility
  - Privacy and security
  - Research using data

Collaborative wiki – dynamic report appendix

# Case study data sources that moved onto acquisition

## Education – State Longitudinal Data Systems

- Washington
- North Carolina
- Virginia
- Kentucky
- Texas



Data Availability
- Non-operational
- Functional
- Student-level available
- One linkage
- Linkage for K–12, Higher Ed, and Workforce

## Housing

- Local Count Real Estate Assessment Data
- Black Knight Financial Services
- CoreLogic
- MLS Data
- Local Crime Reports
- Location Inc.

# Focus Briefly on Data Profiling

- Determination of quality of data and utility <span style="color:darkred">to project at hand</span>
  - Data Quality
  - Data Structure
  - Metadata and Provenance
- Issues <span style="color:darkred">are discovered but not fixed</span>, fix depends on research need

# Local Data: Education Example

# Exercising the data framework for the education data

- ## Profiled variables
  - Student ID, district code, year, gender, race/ethnicity, grade, age, and other variables (e.g., limited English proficiency)
  - Most variables were valid and consistent requiring very little cleaning

- ## Transformed variables
  - Matched school districts with counties
  - Calculated ages from birthdates (NC and KY)
  - Texas enrollment estimates weighted to match the state level counts

- ## Restructured variables
  - Virginia data were restructured to create three main tables for race/ethnicity, by grade, gender by grade, and disadvantaged status by grade

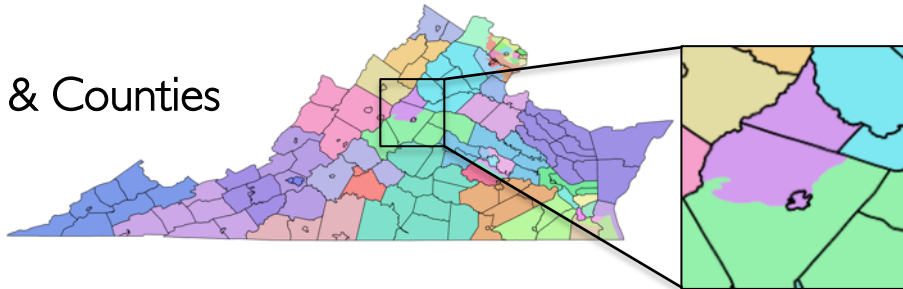# Local data challenges with geographic alignment



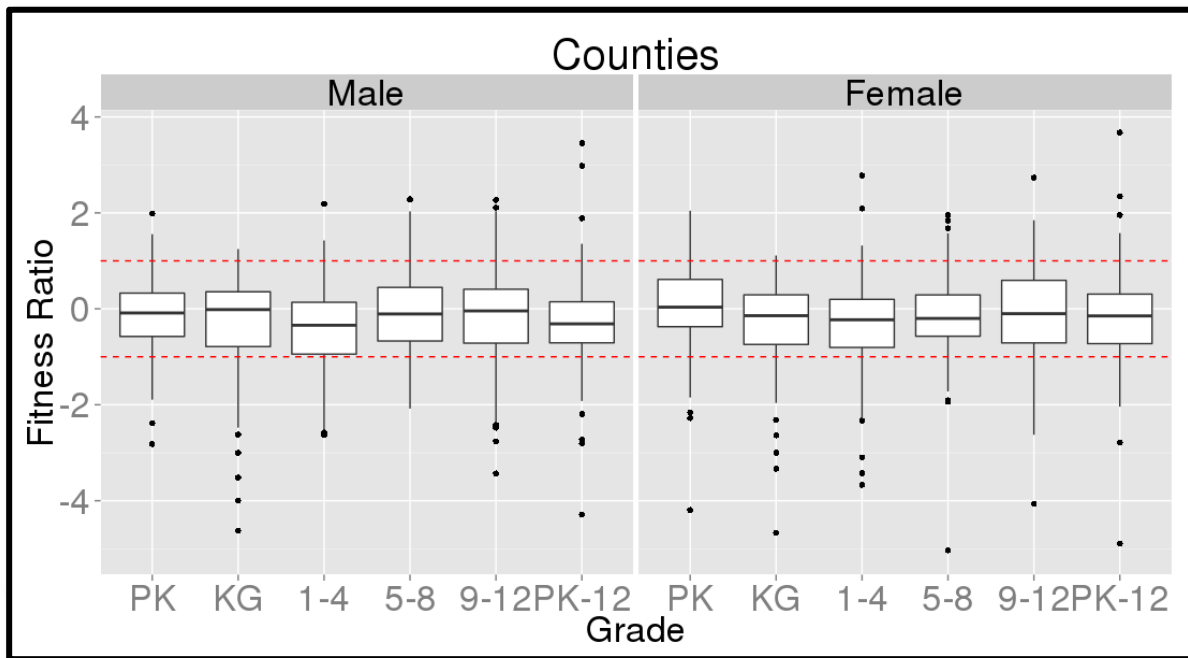Counties & School districts

PUMAs & Counties

PUMAs & Counties

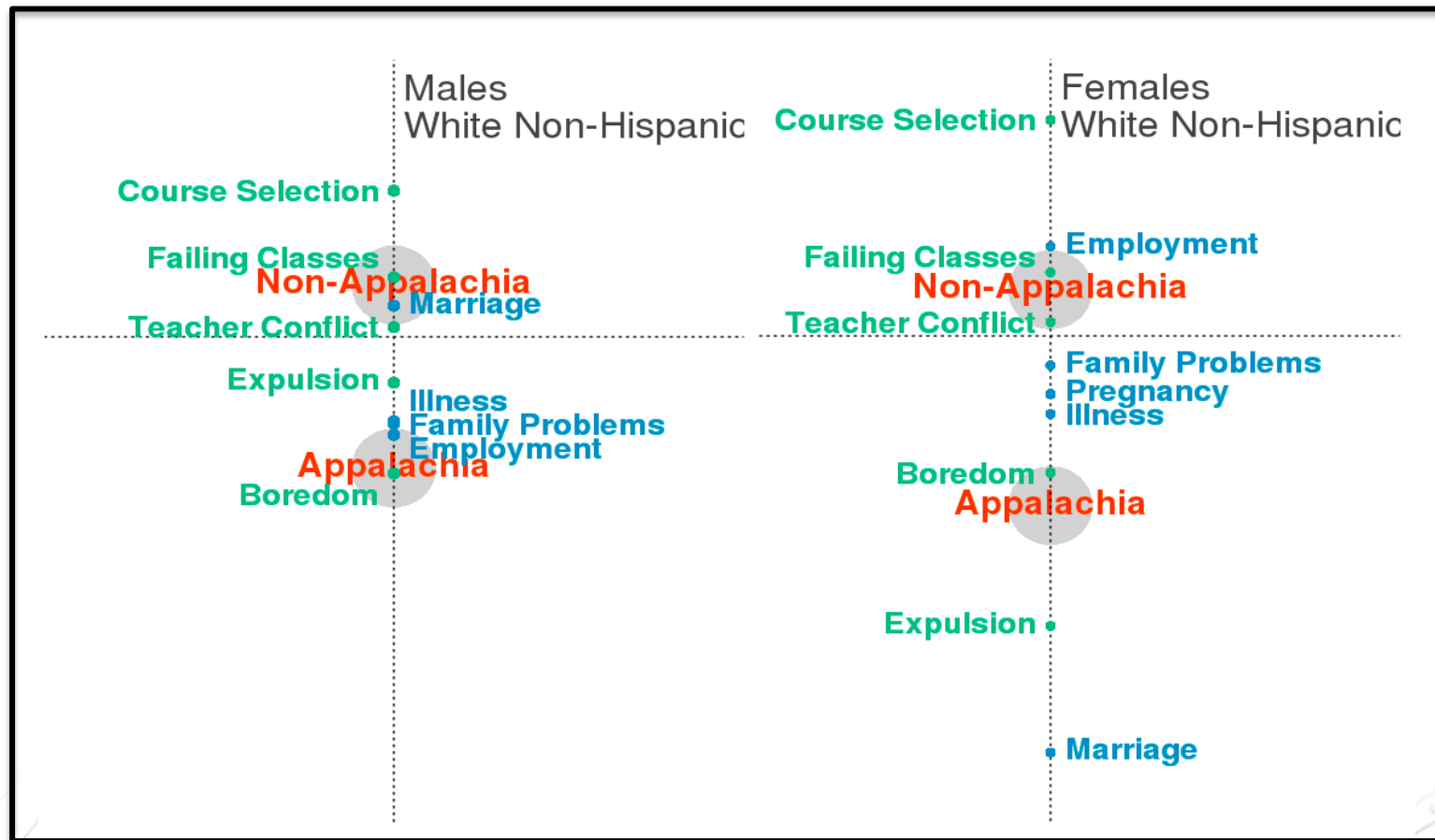# Education SLSD Alignment with ACS – Kentucky

$$\textit{Fitness Ratio} = \frac{\textit{ACS Estimate} - \textit{Local Data Estimate}}{\textit{ACS 90\% Margin of Error}}$$

## Public Enrollment at Kentucky State and County Levels, 2009-2013

# Kentucky Public School Dropouts
## Grades 9-12, KLDS 2009-2013

# Local Data: Housing Example

# Exercising the data framework for the housing data

- **Profiled variables**
  - Need to align units of analysis, parcels versus housing units
  - Reconcile land use codes through discussion with experts
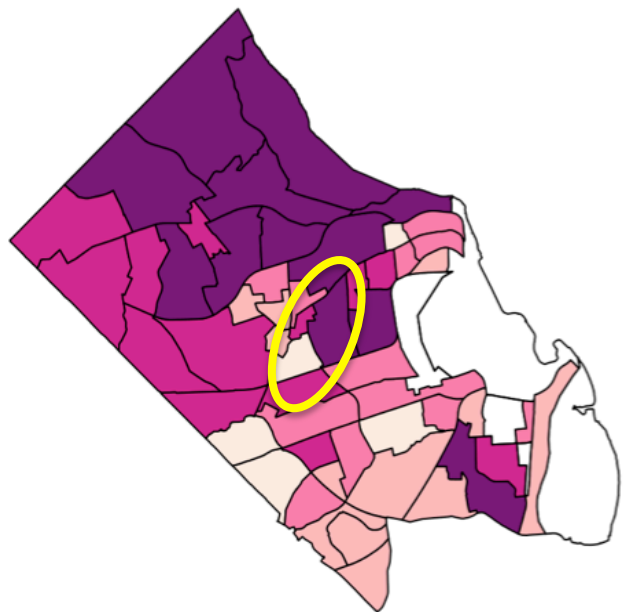  - Remove non-residential parcels

- **Transformed variables**
  - Weight parcels by (estimated) number of units
  - Adjusted inconsistent longitudinal data, e.g., number of units or year built
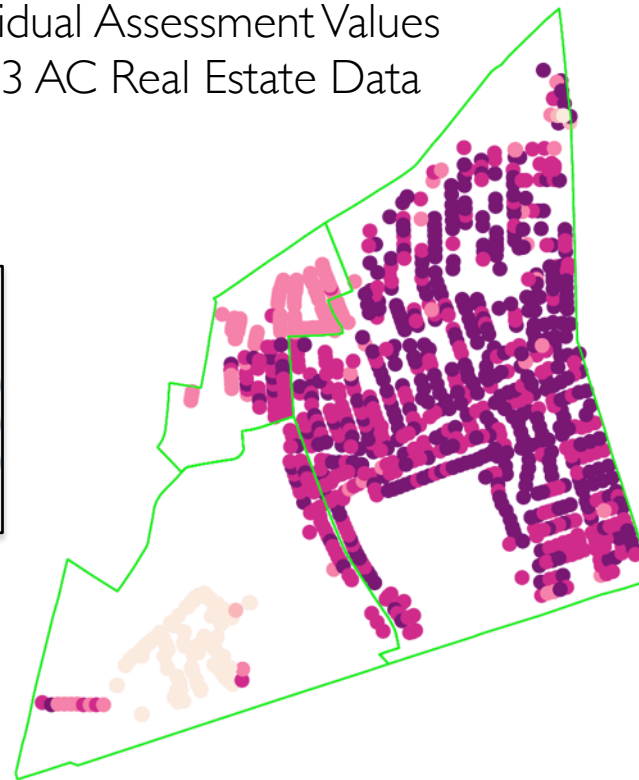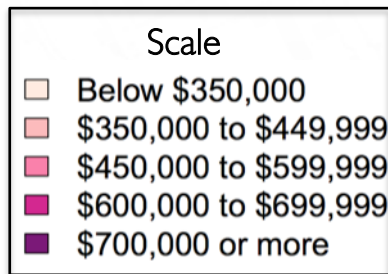
- **Restructured variables**
  - Create a consistent set of geocodes per parcel

# Characterizing neighborhoods

Medan House Value
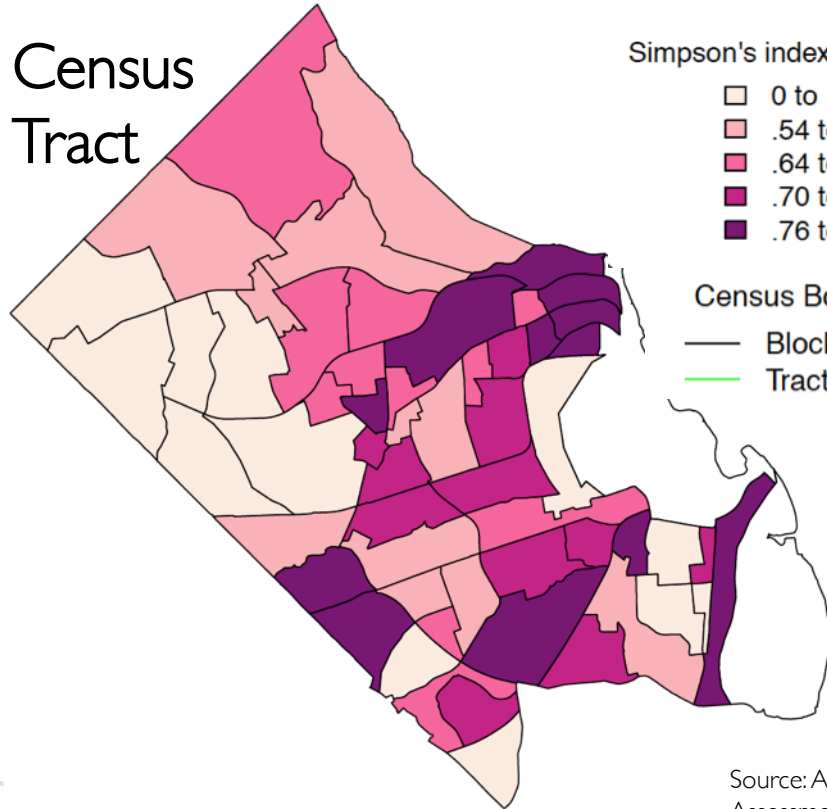Owner Occupied Units
ACS 2009-2013 Estimate

Individual Assessment Values
2013 AC Real Estate Data



### Scale

| | |
|---|---|
| ☐ | Below $350,000 |
| ☐ | $350,000 to $449,999 |
| ☐ | $450,000 to $599,999 |
| ☐ | $600,000 to $699,999 |
| ☐ | $700,000 or more |

Source: Arlington County Real Estate
Assessment, 2013; ACS 2009-2013

# Simpson Index of Home Value Diversity



Census Tract

Census Block Group

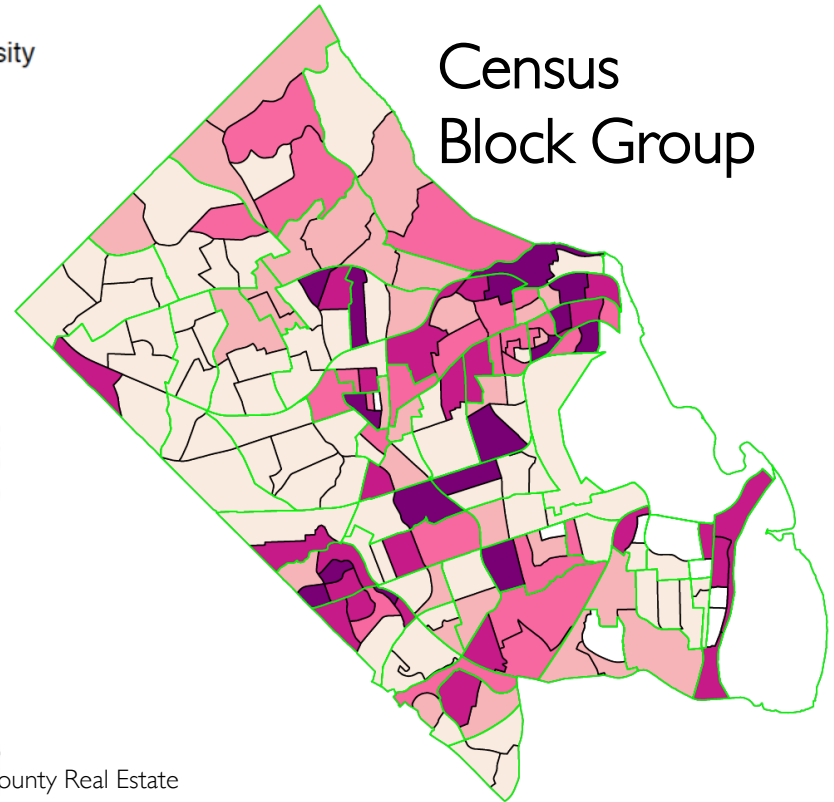**Simpson's index of diversity**
- 0 to .53
- .54 to .63
- .64 to .59
- .70 to .75
- .76 to .86

**Census Borders**
— Block Group
— Tract

Source: Arlington County Real Estate Assessment 2013

# Concluding Remarks

- Use of external data – no control over collection, unlike federal statistics

- Through lens of case studies develop and test data framework

- Need a disciplined, yet flexible and adaptable, data framework to assess data quality and fitness-for-use that is dependent on use, e.g., Innovation measures

# Proposed Criteria for Future Research

- Speed and regularity of acquiring data

- Flexibility in negotiating data agreements

- Capacity to implement Data Framework
  - e.g., Organizational structure and constraints, GIS, longitudinal discrepancies, statistical data integration, etc.

- Research needed to align and harmonize data sources

- Value proposition to sponsor