

# A Kerfuffle

## Differential Privacy and the 2020 Census

Aloni Cohen  
University of Chicago

Workshop on Societal Considerations & Applications

Simons Institute

8 Nov 2022



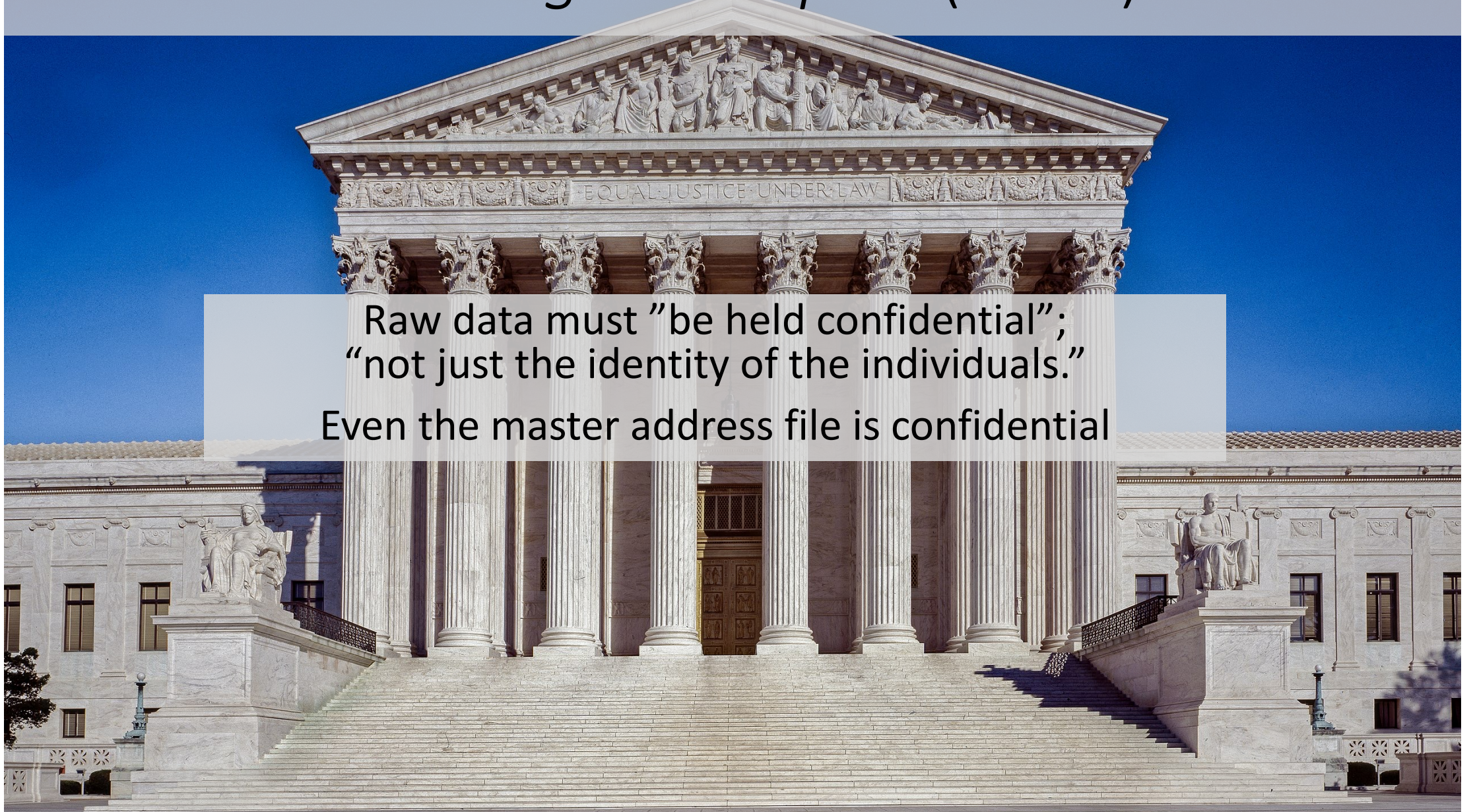
# The Census Act (1976)

The background of the slide is a photograph of the United States Capitol building in Washington, D.C. The building is a grand neoclassical structure with a prominent white dome topped by a statue. The facade is supported by a series of tall, white columns. The sky is a clear, bright blue.

The Census Bureau may not  
“disclose the information reported by, or on  
behalf of, any particular respondent”

# *Baldrige v Shapiro (1982)*

Raw data must "be held confidential";  
"not just the identity of the individuals."  
Even the master address file is confidential



**Implications of Differential Privacy for  
Census Bureau Data and Scientific Research**

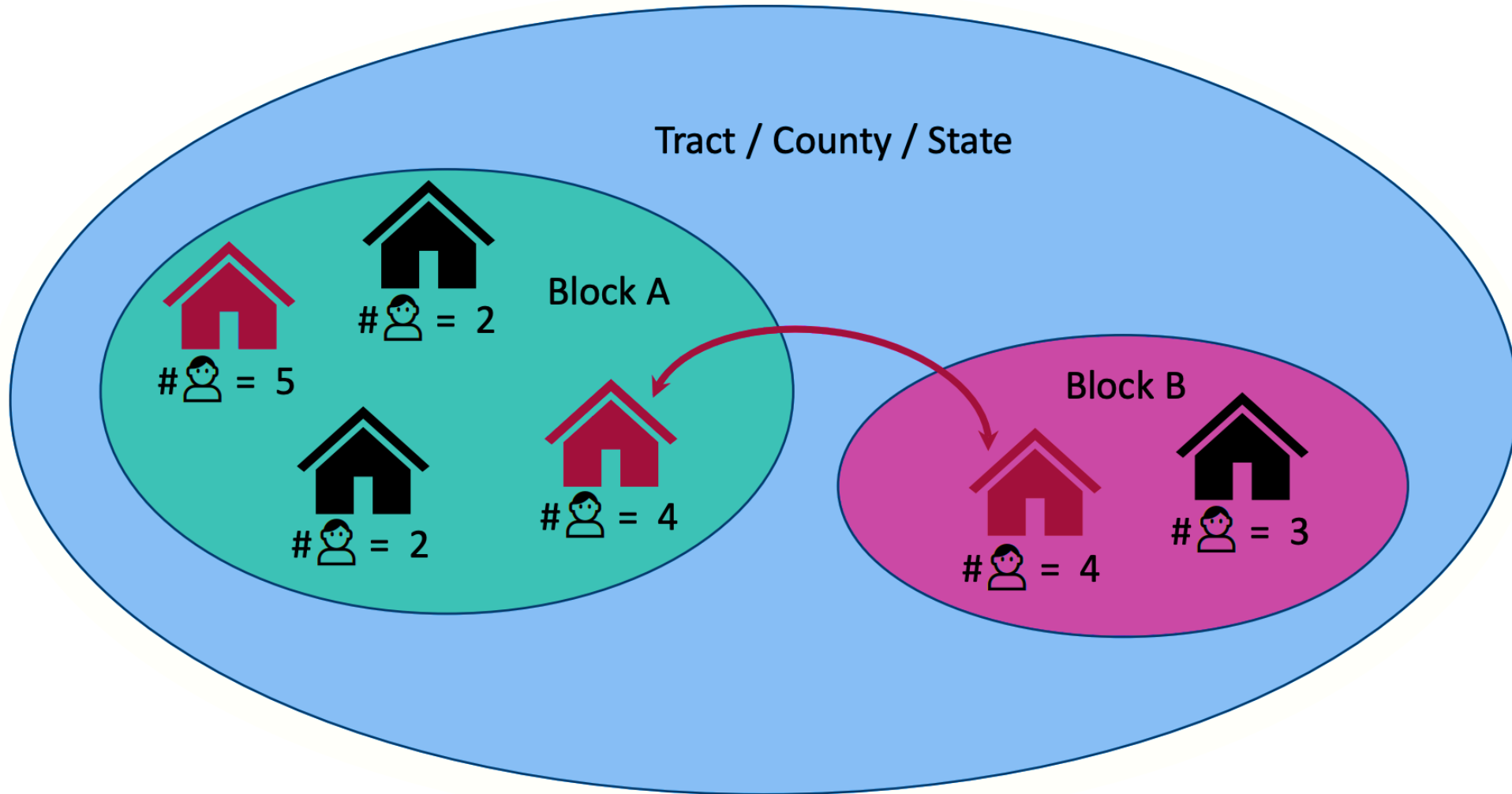
Task Force on Differential Privacy for Census Data†  
Institute for Social Research and Data Innovation (ISRDI)  
University of Minnesota

December 2018  
Version 5.1  
Working Paper No. 2018-6  
<https://doi.org/10.18128/MPC2018-6>

This report was prepared by Steven Ruggles (ISRDI) with the assistance of Margo J. Anderson (University of Wisconsin-Milwaukee), Jane Bambauer (Arizona State University), Michael Davern (NORC), Reynolds Farley (University of Michigan),

The meaning of the law is clear and unambiguous: census publications must ensure that the responses of particular identified persons cannot be determined from census publications. To comply with the law, it is not necessary to mask the *characteristics* of individuals; rather, it is necessary to mask the *identity* of individuals. Thus, for the past six decades the Census Bureau disclosure control strategy has focused on targeted strategies to prevent re-identification attacks, so that an outside adversary cannot positively identify which person provided a particular response. The protections in place—sampling, swapping, suppression of geographic information and extreme values, imputation, and perturbation—have worked extremely well to meet this standard. Indeed, there is not a single documented case of anyone outside the Census Bureau revealing the responses of a particular identified person by breaking into public use decennial census or ACS data.

# Swapping (1990 – 2010)



# 2010 Reconstruction



**John Abowd** @john\_abowd · Apr 7, 2019

Replying to @john\_abowd

1. First, let's get the facts straight: the U.S. Census Bureau reconstructed 100% of the 2010 Census micro-data records (308,745,538 persons).



2



1



6



**John Abowd** @john\_abowd · Apr 7, 2019

3. The reconstructed records matched the confidential data (2010 CEF) exactly (every single bit) for 46% of the population (142 million people) and allowing age +/- 1 year for 71% of the population (219 million people).



1



5



**Steven Ruggles** @HistDem · Apr 20, 2021

1. I prepared a report for the Plaintiffs in the Alabama v. Department of Commerce lawsuit over differential privacy in the census, available here: [users.hist.umn.edu/~ruggles/censi...](https://users.hist.umn.edu/~ruggles/censi...)

**UNITED STATES DISTRICT COURT FOR THE  
MIDDLE DISTRICT OF ALABAMA  
EASTERN DIVISION**

THE STATE OF ALABAMA, *et al.*,

Plaintiffs,

v.

UNITED STATES DEPARTMENT OF  
COMMERCE; GINA RAIMONDO, *et al.*,

Defendants.

CASE NO. 3:21-cv-00211-RAH-ECM-KCN



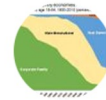
1



16



17



**Steven Ruggles** @HistDem · Apr 20, 2021

2. I argue that the database reconstruction experiment did not demonstrate a convincing threat to confidentiality, because the results reported by the Census Bureau can be largely explained by chance.



1



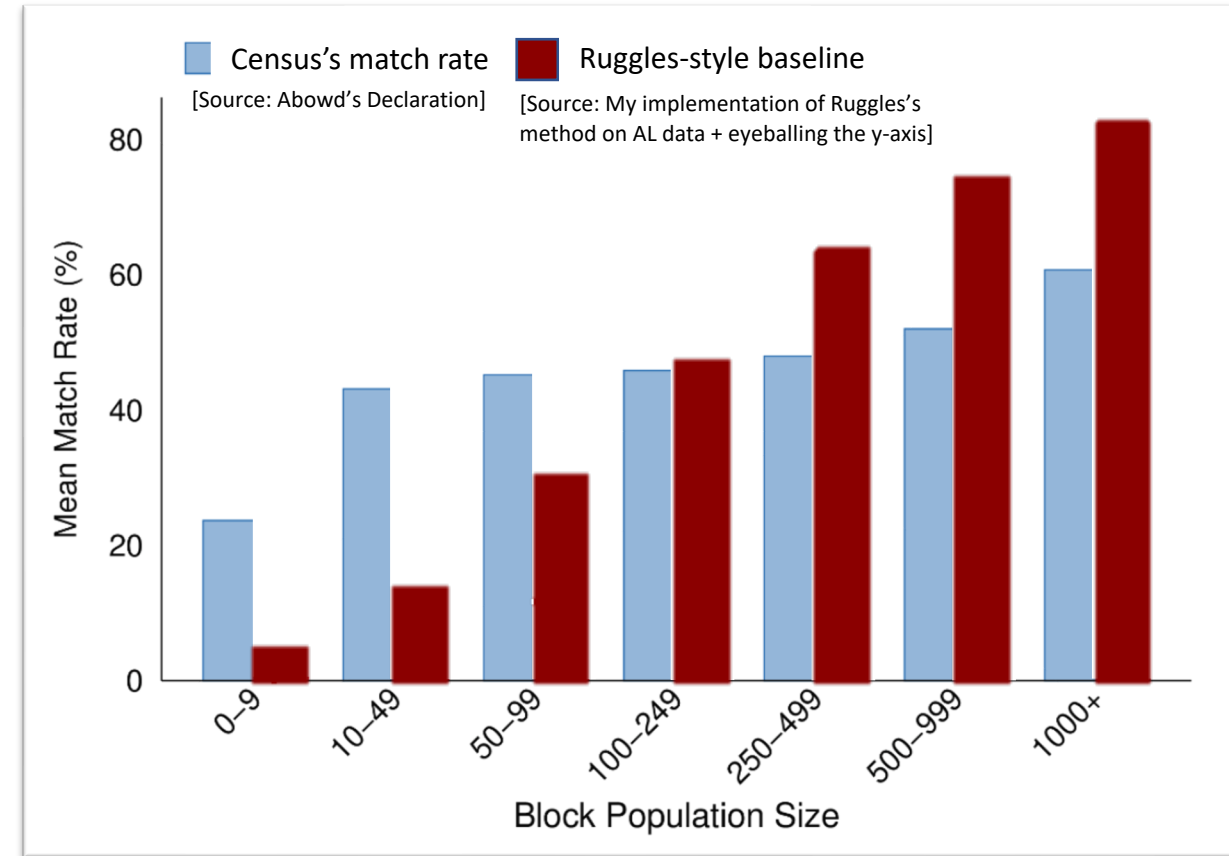
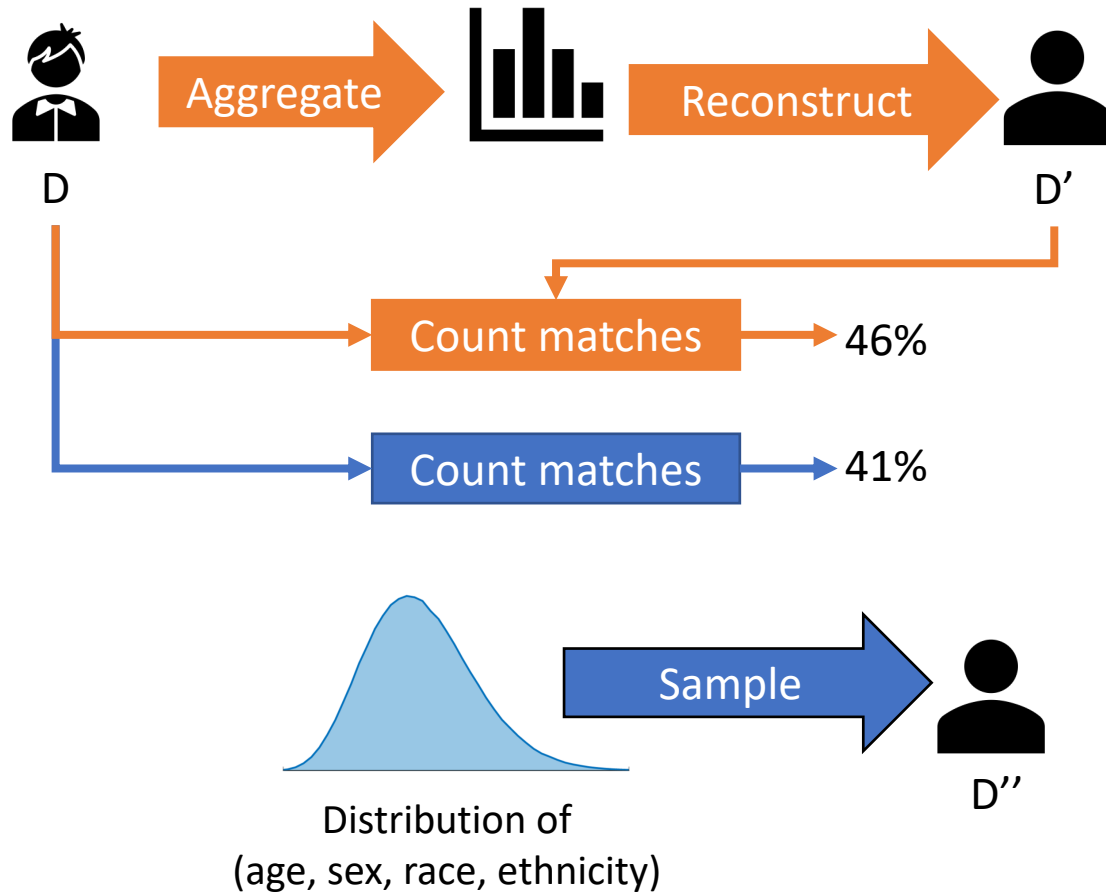
1



2



# 2010 Reconstruction



From: Filings in Alabama v Dept of Commerce; Ruggles, Van Riper, "The Role of Chance in the Census Bureau Database Reconstruction Experiment"

# 2020 Census Operational Plan

---

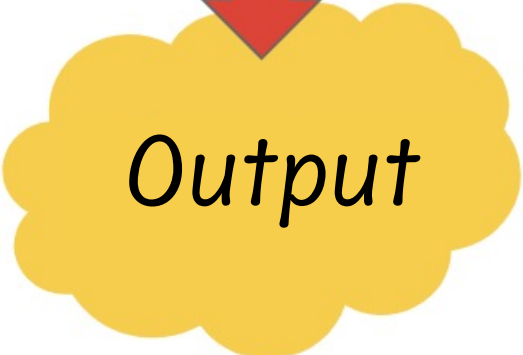
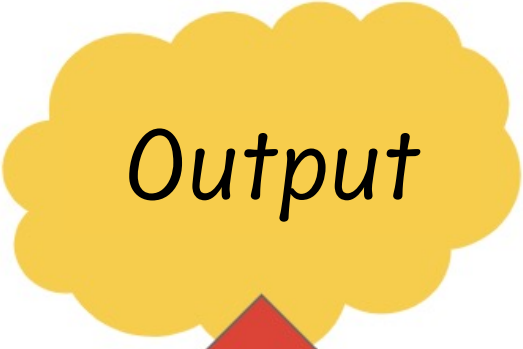
*A New Design for the 21st Century*

Issued January 2022  
Version 5.0

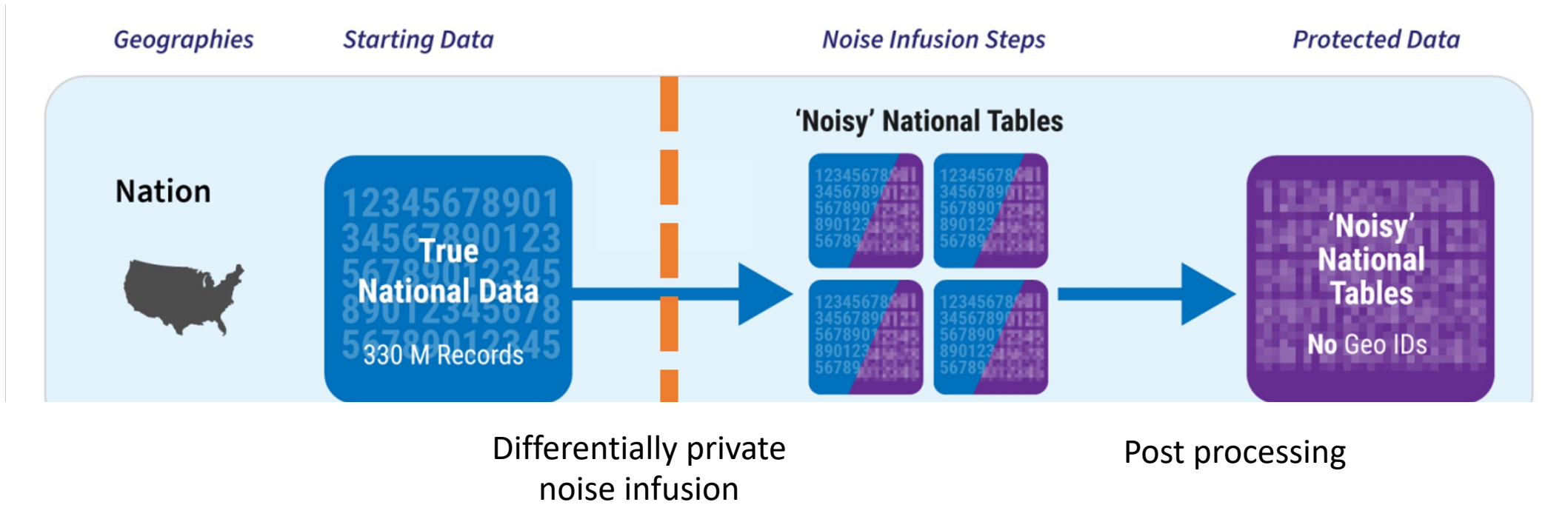




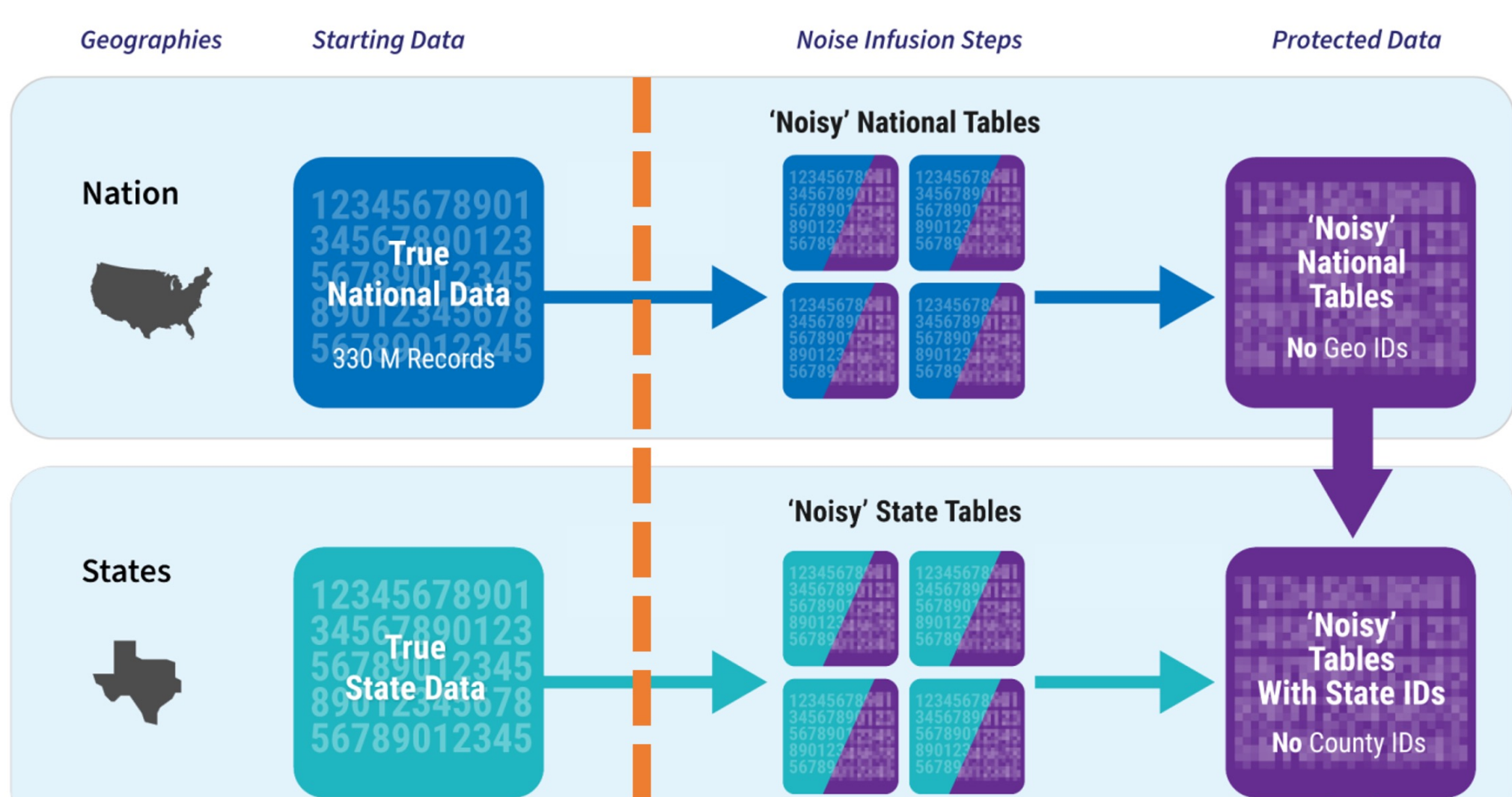
# Differential privacy



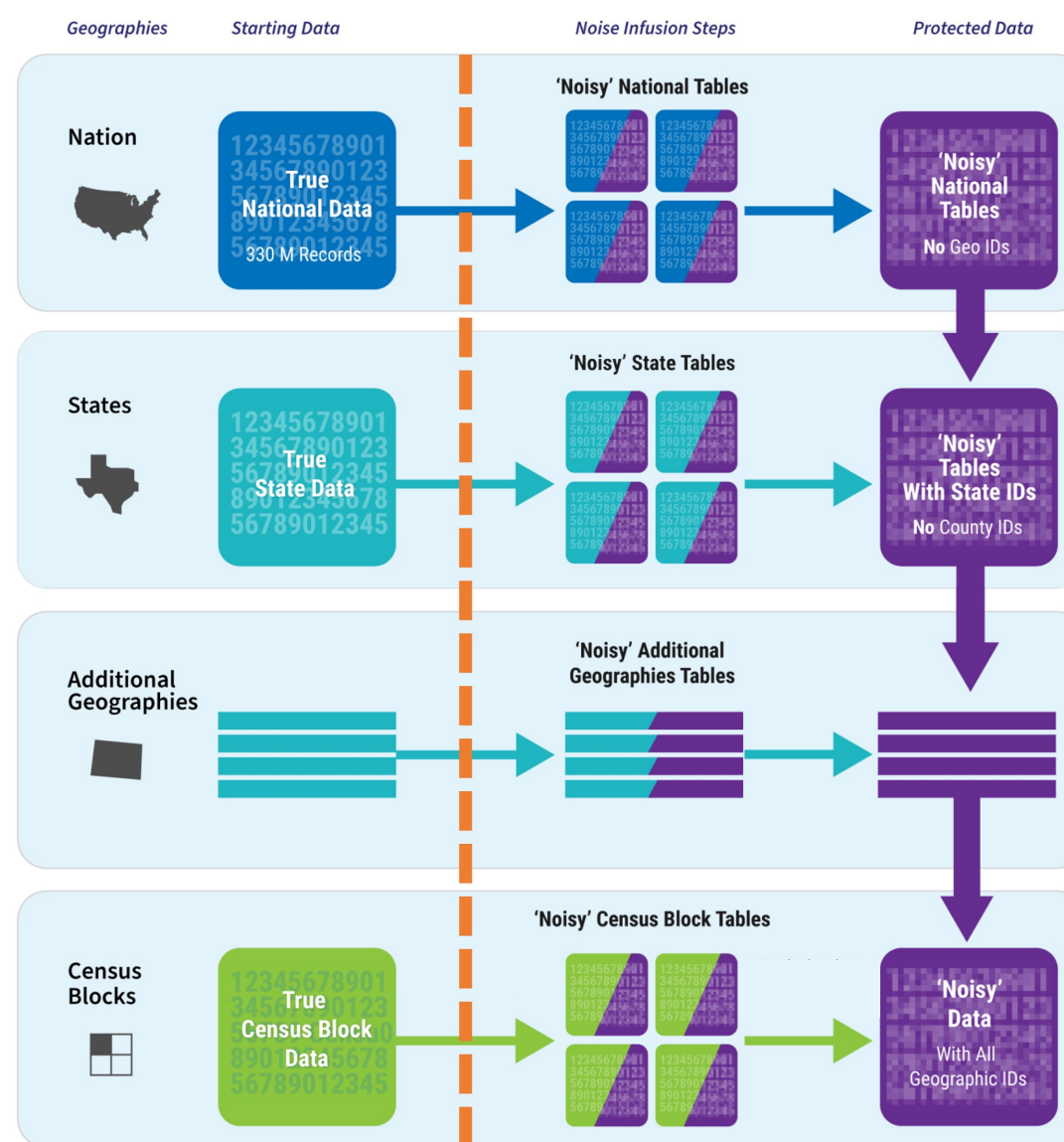
# TopDown Algorithm



# TopDown Algorithm



# TopDown Algorithm



**Table 2.3. Hypothetical Example of Post-Processing**

**Step 2: Post-processing**

| Block              | Enumerated counts                  |   |                          | Noise                              |   |                          | Preliminary noisy counts           |   |                          | Post-processed counts              |   |                          |
|--------------------|------------------------------------|---|--------------------------|------------------------------------|---|--------------------------|------------------------------------|---|--------------------------|------------------------------------|---|--------------------------|
|                    | Popu-<br>lation<br>under<br>age 18 | Popu-<br>lation<br>aged<br>18 and<br>over | Total<br>popu-<br>lation | Popu-<br>lation<br>under<br>age 18 | Popu-<br>lation<br>aged<br>18 and<br>over | Total<br>popu-<br>lation | Popu-<br>lation<br>under<br>age 18 | Popu-<br>lation<br>aged<br>18 and<br>over | Total<br>popu-<br>lation | Popu-<br>lation<br>under<br>age 18 | Popu-<br>lation<br>aged<br>18 and<br>over | Total<br>popu-<br>lation |
| Block 1.....       | 25                                 | 75  | 100                      | 0                                  | -4  | 2                        | 25                                 | 71  | 102                      | 27 (+2)                            | 71 (-4)                                   | 98 (-2)                  |
| Block 2.....       | 20                                 | 70  | 90                       | -3                                 | 2   | 3                        | 17                                 | 72  | 93                       | 19 (-1)                            | 72 (+2)                                   | 91 (+1)                  |
| Block 3.....       | 10                                 | 40  | 50                       | 2                                  | -3  | -2                       | 12                                 | 37  | 48                       | 12 (+2)                            | 37 (-3)                                   | 49 (-1)                  |
| Block 4.....       | 1                                  | 9   | 10                       | -2                                 | 1   | 1                        | -1                                 | 10  | 11                       | 0 (-1)                             | 11 (+2)                                   | 11 (+1)                  |
| Block 5.....       | 1                                  | 2   | 3                        | 0                                  | 2   | 0                        | 1                                  | 4   | 3                        | 1 (+0)                             | 4 (+2)                                    | 5 (+2)                   |
| <b>Block group</b> |                                    |   |                          |                                    |   |                          |                                    |   |                          | <b>59</b>                          | <b>195</b>                                | <b>254</b>               |

Source: U.S. Census Bureau.

# Noise from TopDown is relatively small

| <b>Error Statistics for Total Population for Counties (Excluding Puerto Rico)</b> |                    |  |   |           |
|---|--------------------|--|---|-----------|
| Counties by size  | Number of counties | Mean absolute error (counts of people) | Error: middle 90 percent (counts of people) |           |
|   |                    |  | Minus                                       | Plus      |
| <b>All counties. ....</b>   | <b>3,143</b>       | <b>1.75</b>                            | <b>-4</b>                                   | <b>+4</b> |
| Counties with total population between 0-249 .....                                | 2                  | 2.00                                   | -1  | +3        |
| Counties with total population between 250-749 .....                              | 19                 | 1.32                                   | -2  | +2        |
| Counties with total population between 750-1,249.....                             | 26                 | 1.38                                   | -2  | +4        |
| Counties with total population between 1,250-1,749. ...                           | 24                 | 1.00                                   | -2  | +3        |
| Counties with total population between 1,750-1,949 ...                            | 14                 | 1.14                                   | -1  | +2        |
| Counties with total population between 1,950-2,049. ...                           | 10                 | 1.50                                   | -1  | +5        |
| Counties with total population between 2,050-2,249 ..                             | 16                 | 0.88                                   | -1  | +1        |
| Counties with total population between 2,250-2,749. ...                           | 35                 | 1.31                                   | -2  | +3        |
| Counties with total population between 2,750-3,249. ...                           | 38                 | 1.29                                   | -2  | +3        |
| Counties with total population at or above 3,250 .....                            | 2,959              | 1.79                                   | -4  | +4        |

### Coverage Error

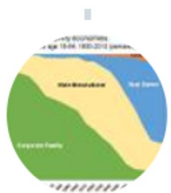
|                           |              |               |               |               |
|---------------------------|--------------|---------------|---------------|---------------|
| <b>All counties. ....</b> | <b>3,143</b> | <b>964.00</b> | <b>-1,841</b> | <b>+2,048</b> |
|---------------------------|--------------|---------------|---------------|---------------|

### Nonsampling Variability

|                           |              |               |             |             |
|---------------------------|--------------|---------------|-------------|-------------|
| <b>All counties. ....</b> | <b>3,143</b> | <b>117.27</b> | <b>-248</b> | <b>+230</b> |
|---------------------------|--------------|---------------|-------------|-------------|



Apocalypse Now



**Steven Ruggles** @HistDem · Jul 5, 2019



I am increasingly convinced that DP will degrade the quality of data available about the population, and will make scientifically useful public use microdata impossible. 3/



2



12

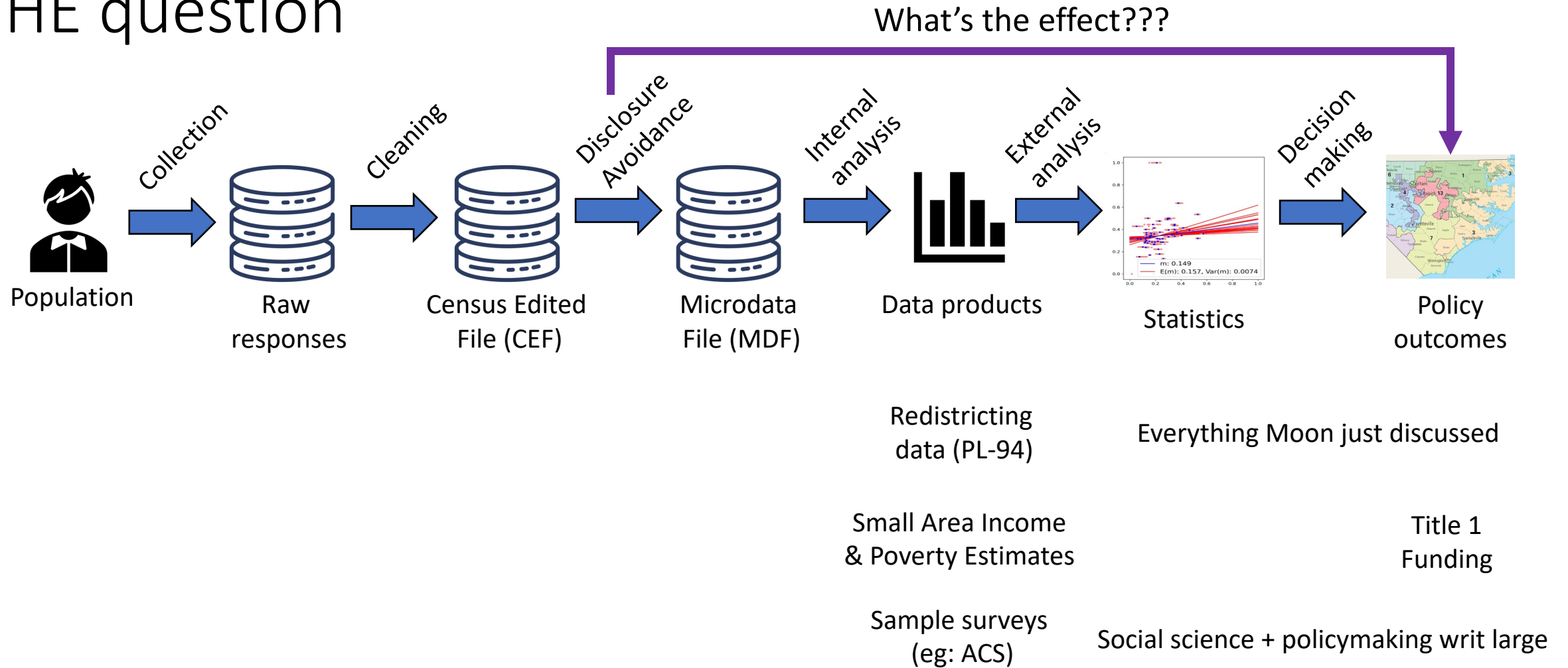


32

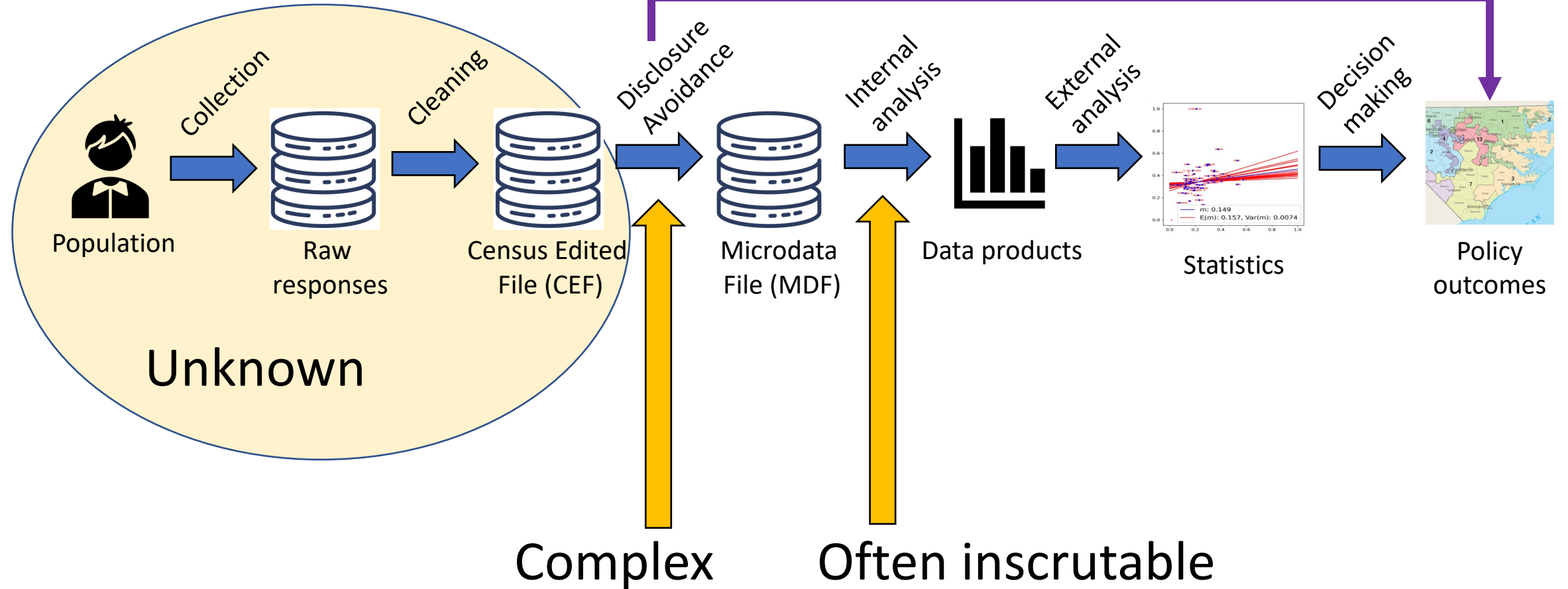




# THE question



# THE question



# The quagmire



# Option 0: Be the Census

**Table 11V.** Counts & Measures of Variation for Tate County School Districts, MS Twenty-five Runs of the *TDA* for County Districts 01, 02, 03, 04, 05  
 ( $C_T(g)$  counts result from 2020 Census Redistricting Data Production Settings ( $\epsilon = 17.14$  for persons) version of *TDA*.)

| DIST-ID                  | (Measures of Variation) |                       |                      |                      |                      |                      |                      |                      |                      |                      |                      |                      |
|--------------------------|-------------------------|-----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|
|                          | Tate Schools            |                       | 01                   |                      | 02                   |                      | 03                   |                      | 04                   |                      | 05                   |                      |
|                          | $\bar{C}_T(g)$          | $C_S(g)$              | $\bar{C}_T(g)$       | $C_S(g)$             | $\bar{C}_T(g)$       | $C_S(g)$             | $\bar{C}_T(g)$       | $C_S(g)$             | $\bar{C}_T(g)$       | $C_S(g)$             | $\bar{C}_T(g)$       | $C_S(g)$             |
|                          | $\sqrt{V(1)_g}$         | $\sqrt{V(2)_g}$       | $\sqrt{V(1)_g}$      | $\sqrt{V(2)_g}$      | $\sqrt{V(1)_g}$      | $\sqrt{V(2)_g}$      | $\sqrt{V(1)_g}$      | $\sqrt{V(2)_g}$      | $\sqrt{V(1)_g}$      | $\sqrt{V(2)_g}$      | $\sqrt{V(1)_g}$      | $\sqrt{V(2)_g}$      |
| Demographic ( <i>g</i> ) | $RV(1)_g$               | $RV(2)_g$             | $RV(1)_g$            | $RV(2)_g$            | $RV(1)_g$            | $RV(2)_g$            | $RV(1)_g$            | $RV(2)_g$            | $RV(1)_g$            | $RV(2)_g$            | $RV(1)_g$            | $RV(2)_g$            |
| TOTAL                    | 18,815<br>18<br>0.001   | 18,823<br>20<br>0.001 | 3,916<br>22<br>0.006 | 3,914<br>22<br>0.006 | 3,885<br>21<br>0.005 | 3,893<br>23<br>0.006 | 3,644<br>20<br>0.006 | 3,665<br>30<br>0.008 | 3,714<br>26<br>0.007 | 3,697<br>31<br>0.008 | 3,657<br>16<br>0.004 | 3,654<br>16<br>0.004 |
| TOTAL18                  | 13,892<br>17<br>0.001   | 13,893<br>17<br>0.001 | 2,776<br>20<br>0.007 | 2,780<br>21<br>0.007 | 2,833<br>19<br>0.007 | 2,826<br>20<br>0.007 | 2,789<br>14<br>0.005 | 2,799<br>17<br>0.006 | 2,766<br>23<br>0.008 | 2,755<br>26<br>0.009 | 2,728<br>13<br>0.005 | 2,733<br>14<br>0.005 |
| TOTALHISP                | 423<br>9<br>0.021       | 399<br>26<br>0.064    | 95<br>6<br>0.066     | 87<br>10<br>0.118    | 64<br>4<br>0.063     | 63<br>4<br>0.066     | 106<br>8<br>0.073    | 110<br>9<br>0.078    | 51<br>6<br>0.119     | 32<br>20<br>0.631    | 106<br>8<br>0.072    | 107<br>8<br>0.071    |
| TOTALNH                  | 18,392<br>18<br>0.001   | 18,424<br>37<br>0.002 | 3,821<br>22<br>0.006 | 3,827<br>23<br>0.006 | 3,821<br>21<br>0.005 | 3,830<br>23<br>0.006 | 3,537<br>19<br>0.005 | 3,555<br>26<br>0.007 | 3,663<br>24<br>0.007 | 3,665<br>24<br>0.007 | 3,551<br>18<br>0.005 | 3,547<br>18<br>0.005 |
| WHITENH                  | 12,805<br>13<br>0.001   | 12,841<br>39<br>0.003 | 3,387<br>14<br>0.004 | 3,378<br>17<br>0.005 | 1,613<br>14<br>0.009 | 1,628<br>21<br>0.013 | 2,833<br>14<br>0.005 | 2,860<br>30<br>0.011 | 2,276<br>20<br>0.009 | 2,293<br>26<br>0.011 | 2,696<br>12<br>0.005 | 2,682<br>19<br>0.007 |
| BLACKNH                  | 5,394<br>11<br>0.002    | 5,389<br>12<br>0.002  | 373<br>12<br>0.033   | 400<br>30<br>0.074   | 2,158<br>10<br>0.004 | 2,139<br>21<br>0.010 | 678<br>11<br>0.016   | 666<br>16<br>0.024   | 1,363<br>14<br>0.010 | 1,349<br>20<br>0.015 | 822<br>15<br>0.018   | 835<br>20<br>0.024   |

# Option 1: A wrong, but useful model

- Noise  $Z \sim N\left(\frac{1}{\rho^*}\right)$  added to counts in each county / tract / block
- Non-negativity + consistency
  - Small counts biased upwards
  - Large counts biased downwards

$$\mathbb{E}(|Z_h|) \approx \frac{0.8}{\sqrt{\rho}}$$

$$\rho_{pop, county} = 0.213 \rightarrow 1.73$$

$$\rho_{pop, tract} = 0.164 \rightarrow 1.97$$

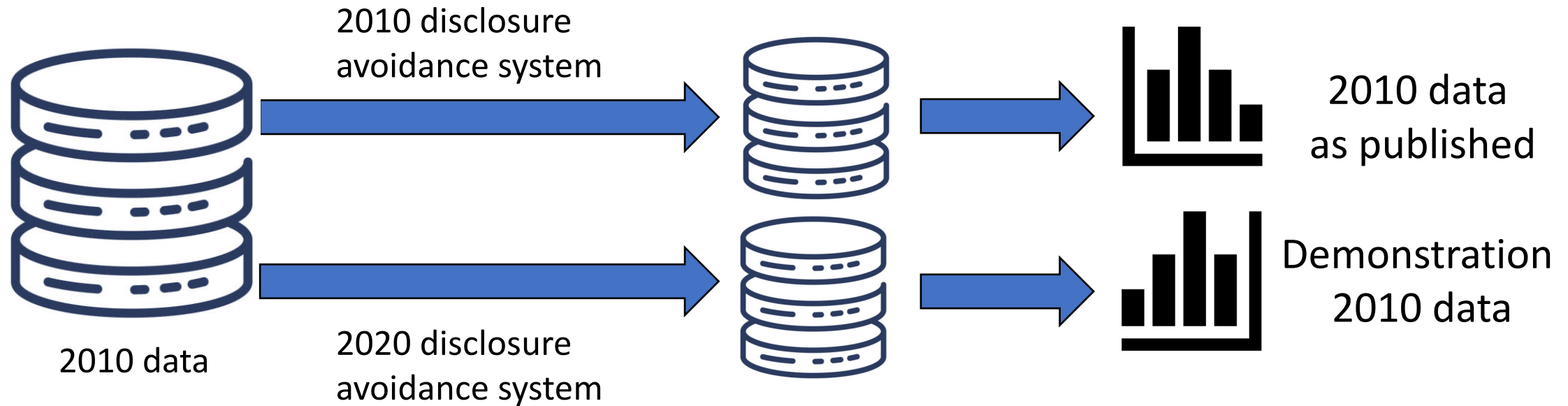
**Privacy-Loss Budget: People**

| Geographic level                         | Rho allocation |
|--|----------------|
| United States .....                      | 104/4,099      |
| State .....                              | 1,440/4,099    |
| County .....                             | 447/4,099      |
| Tract .....                              | 687/4,099      |
| Optimized block group <sup>1</sup> ..... | 1,256/4,099    |
| Block .....                              | 165/4,099      |

| B              | C     | D         | E         | F         | G                 |
|----------------|-------|-----------|-----------|-----------|-------------------|
| name           | state | TOTPOP_dp | TOTPOP_sf | Deviation | Mean( Deviation ) |
| Autauga County | 1     | 54574     | 54571     | 3         | 1.741695126       |
| Baldwin County | 1     | 182266    | 182265    | 1         |                   |
| Barbour County | 1     | 27456     | 27457     | 1         |                   |
| Bibb County    | 1     | 22917     | 22915     | 2         |                   |
| Blount County  | 1     | 57322     | 57322     | 0         |                   |
| Bullock County | 1     | 10915     | 10914     | 1         |                   |

| B                | C     | D         | E         | F         | G                 |
|------------------|-------|-----------|-----------|-----------|-------------------|
| name             | state | TOTPOP_dp | TOTPOP_sf | Deviation | Mean( Deviation ) |
| Census Tract 201 | 1     | 1910      | 1912      | 2         | 1.924650685       |
| Census Tract 202 | 1     | 2171      | 2170      | 1         |                   |
| Census Tract 203 | 1     | 3371      | 3373      | 2         |                   |
| Census Tract 204 | 1     | 4384      | 4386      | 2         |                   |

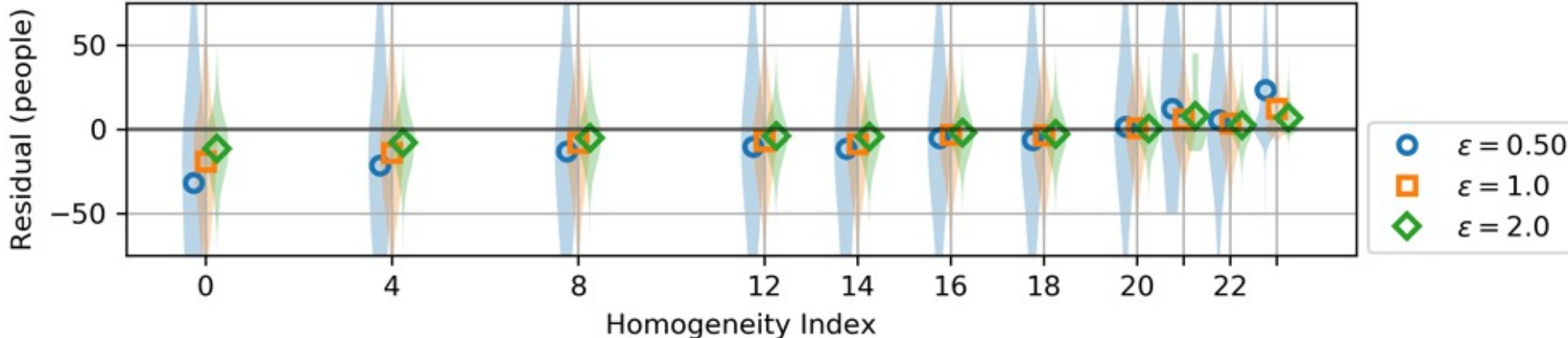
# Option 2: Demonstration data



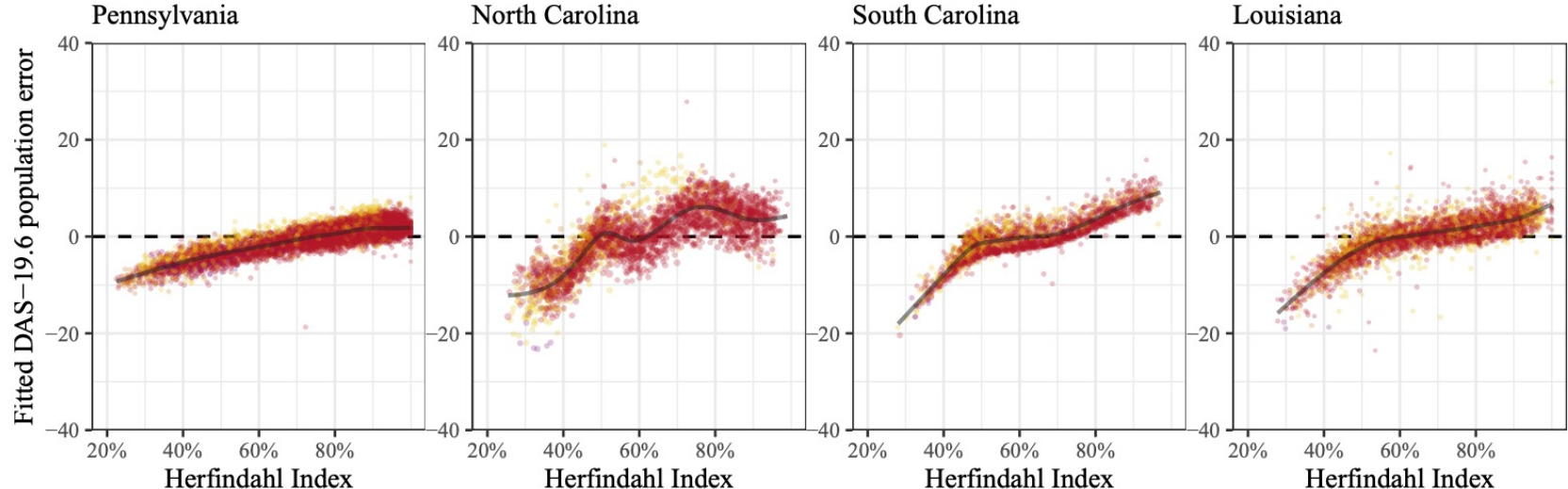
- Descriptive statistics
  - Learn about (TopDown – Swapping)
- Redo prior analysis and compare
  - Learn about effect of noise generally

# Bias from non-negativity

1940 data  
old TopDown

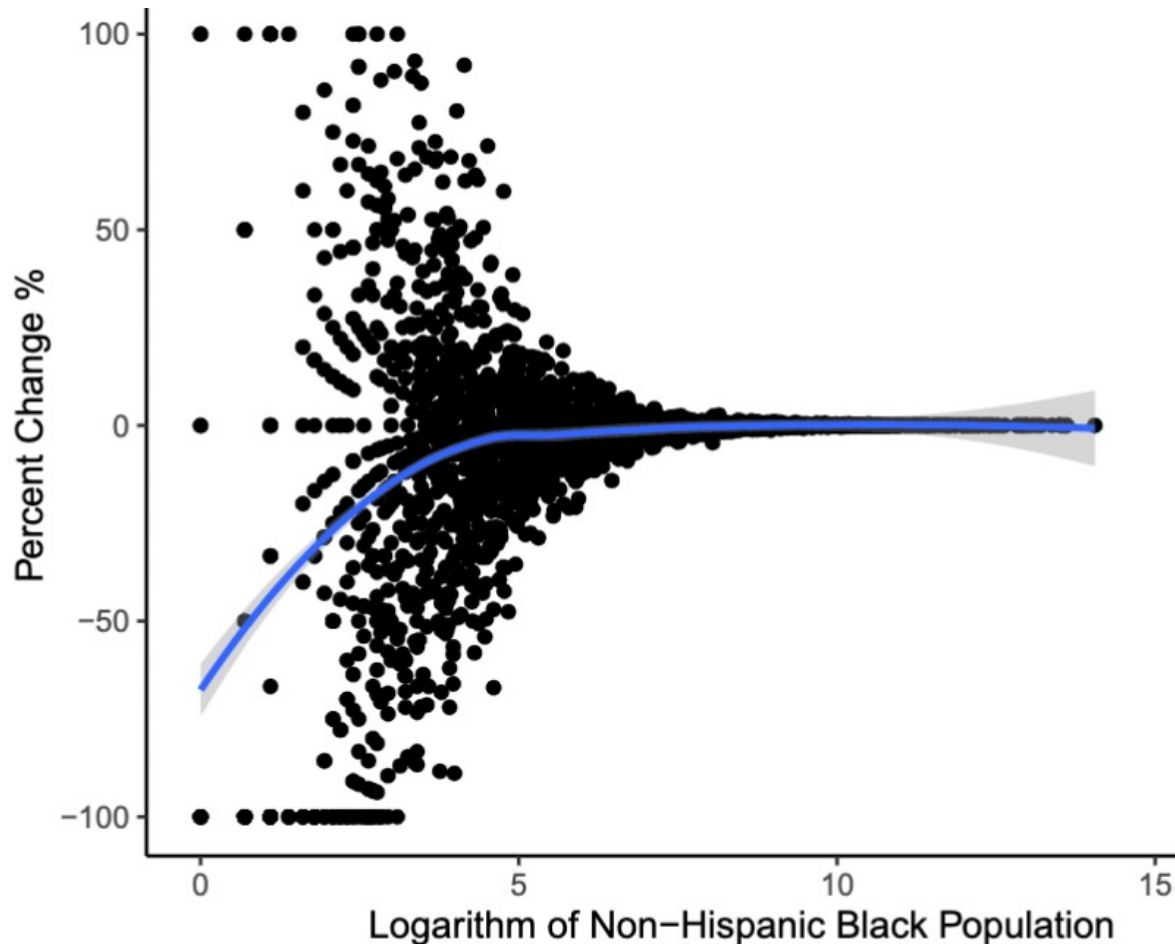


2010  
production  
TopDown



From: Petti, Flaxman "Differential privacy in the 2020 US census: what will it do? Quantifying the accuracy/privacy tradeoff"  
Kenny, Kuriwaki, McCartan, Rosenman, Simko, Imai. "The Use of Differential Privacy for Census Data and its Impact on Redistricting: The Case of the 2020 U.S. Census."

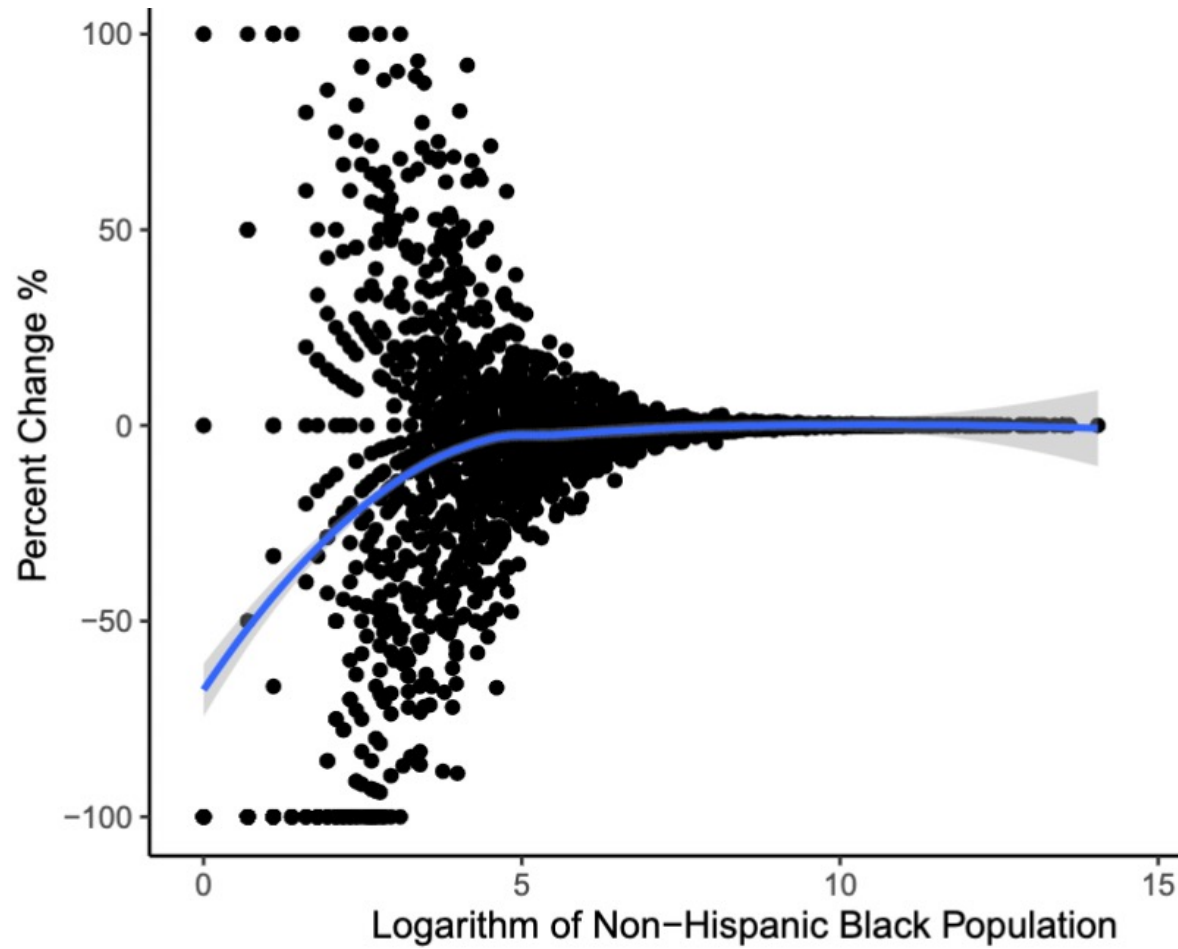
# Mortality rates



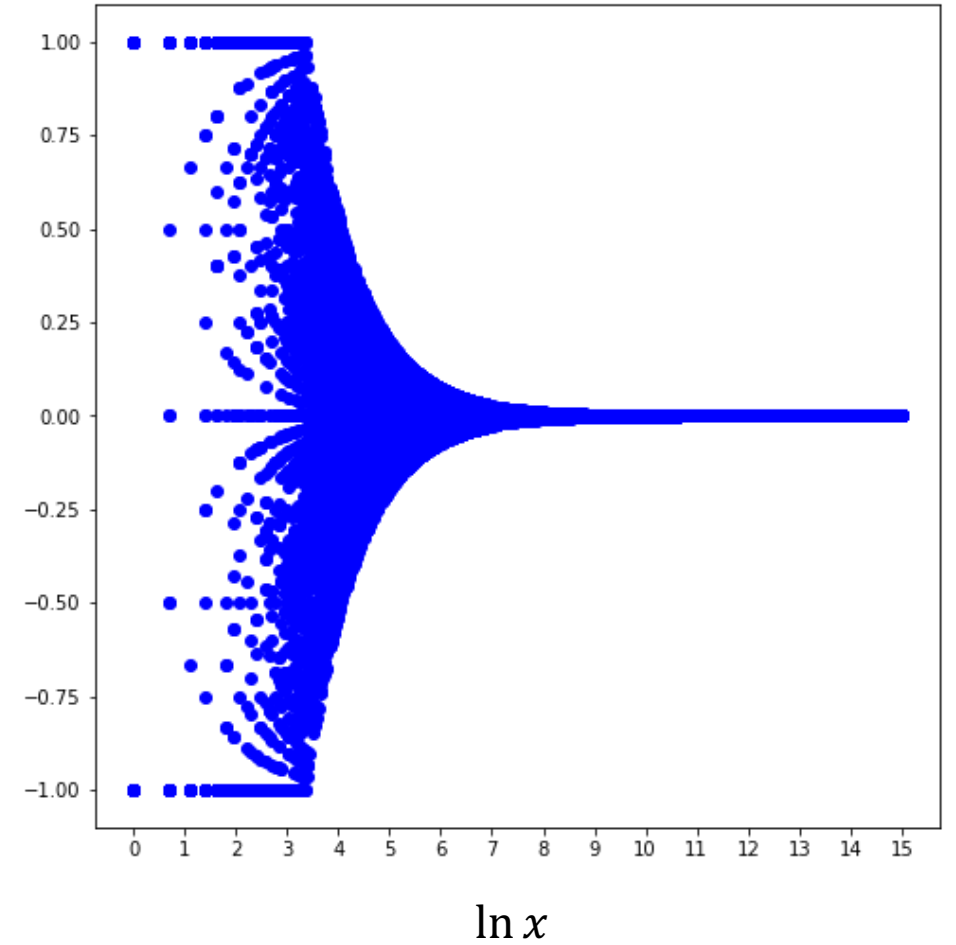
We find that the implementation of differential privacy will produce **dramatic changes in population counts for racial/ethnic minorities in small areas** and less urban settings, significantly altering knowledge about health disparities in mortality.



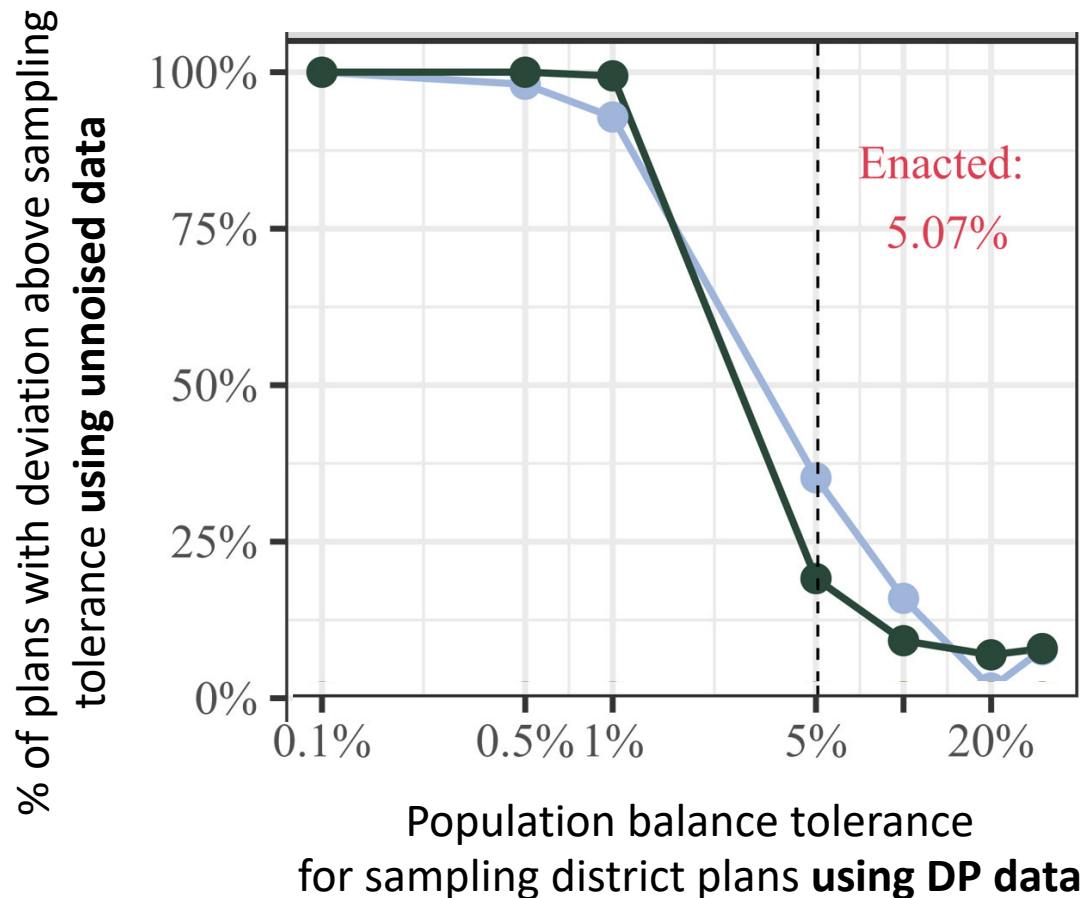
# Mortality rates



$$\frac{Unif(\pm 30)}{x}$$



# Population-balanced redistricting

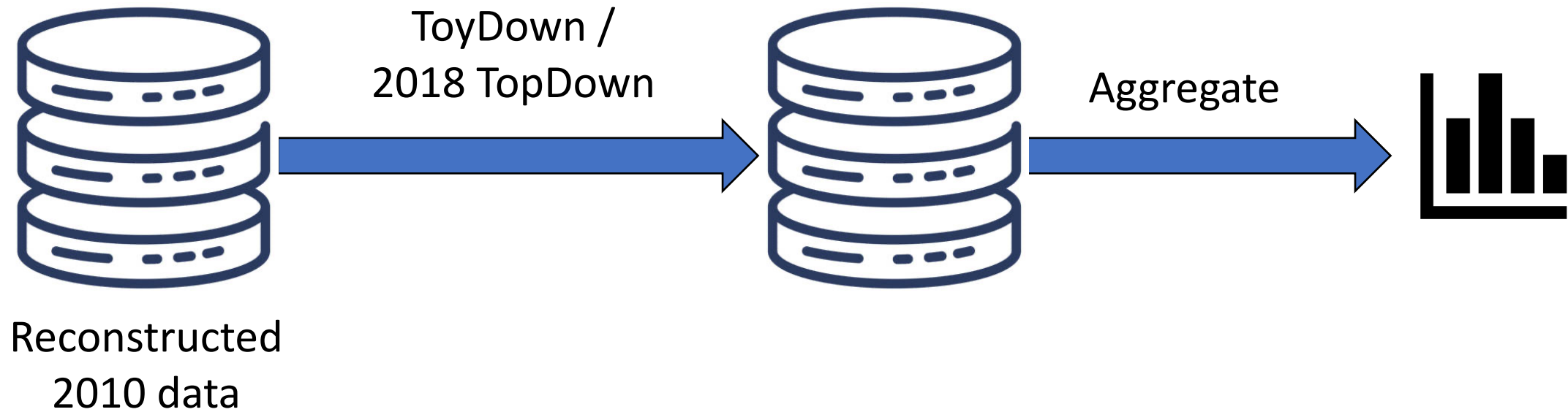


“We conduct our empirical evaluation under a likely scenario, in which practitioners, map drawers and **analysts** alike, **treat these DAS-protected data ‘as is’** as they have done in the past, **without accounting for the DAS noise** generation mechanism.

“Our analysis shows that the added noise makes it **impossible to follow the principle of One Person, One Vote,**”

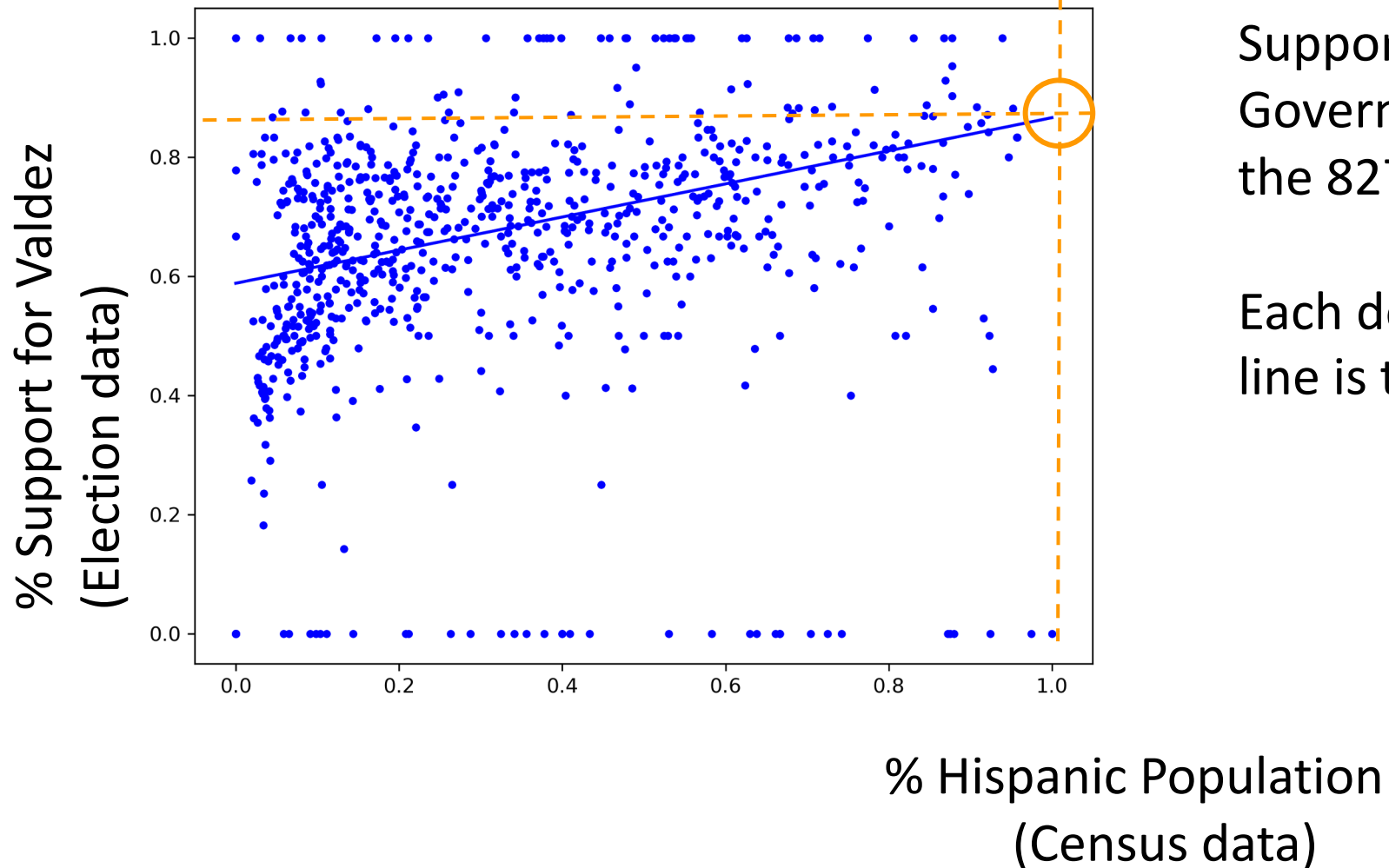
**Solution:** set **sampling** tolerance  $\approx 0.3\%$  lower than **policy** tolerance

# Option 3: Try to run TopDown



# Measuring racial polarized voting (RPV)

Standard ecological regression

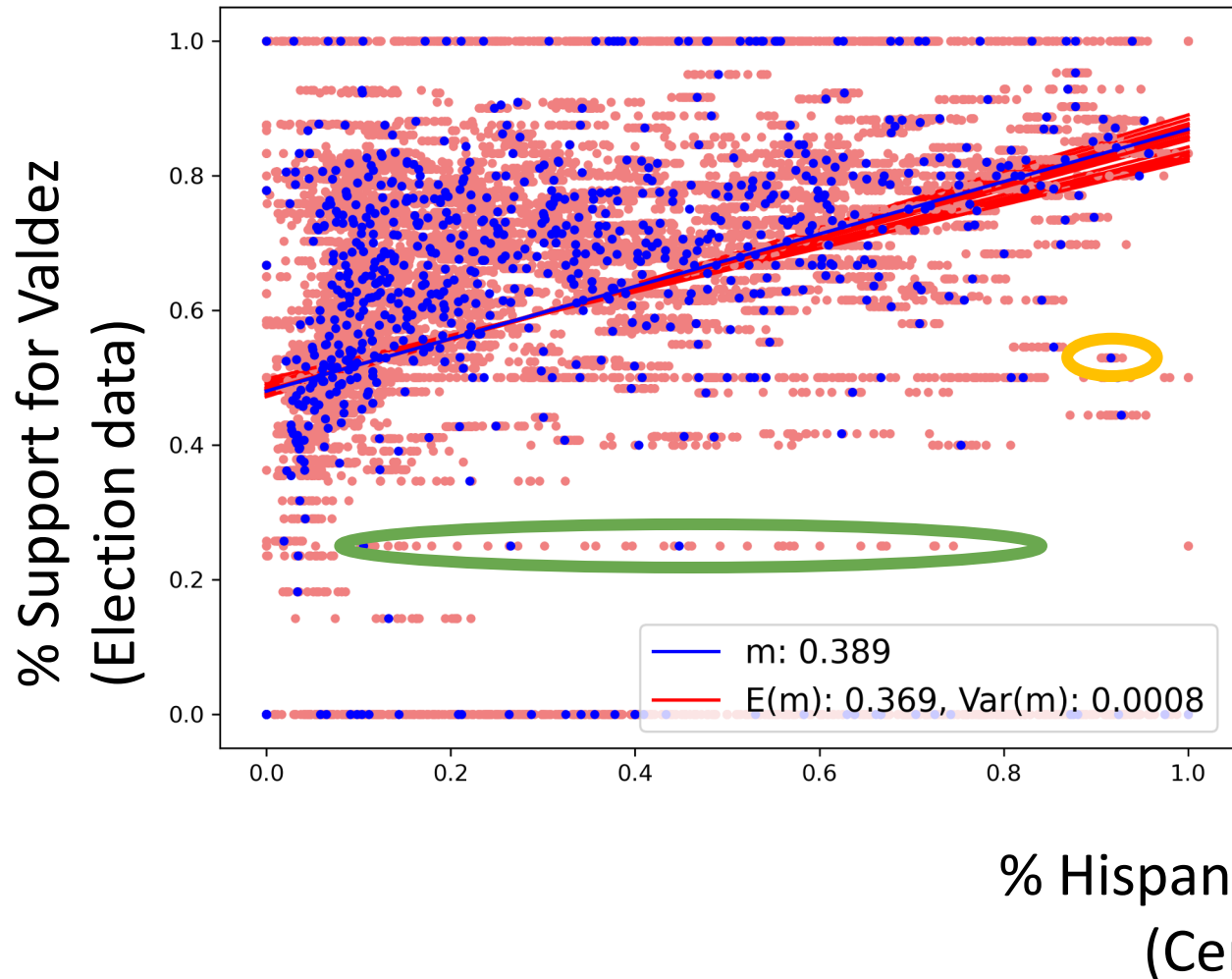


Support for Lupe Valdez in the 2018 Governor election in Texas across the 827 precincts in Dallas County.

Each dot is a precinct, and the blue line is the line of best fit.

# RPV with noise

Standard ecological regression



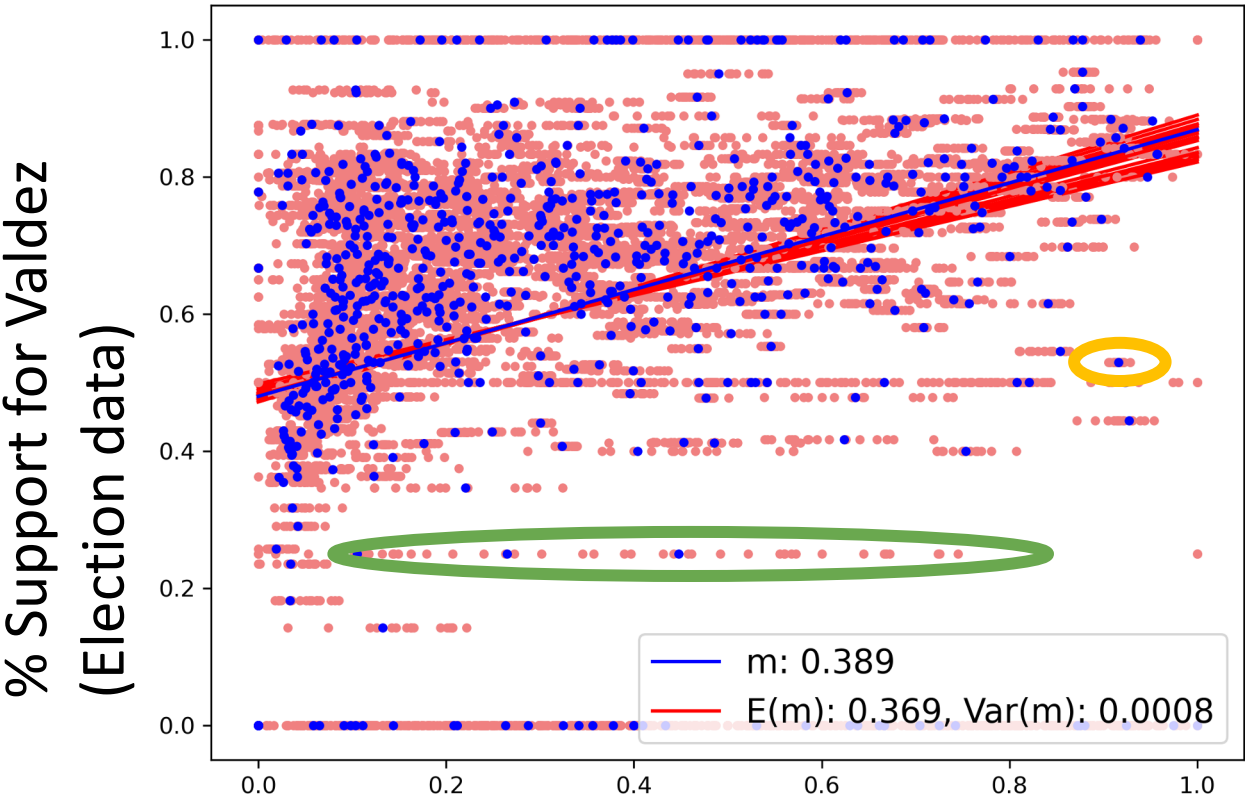
Support for Lupe Valdez in the 2018 Governor election in Texas across the 827 precincts in Dallas County.

Each dot is a precinct, and the blue line is the line of best fit.

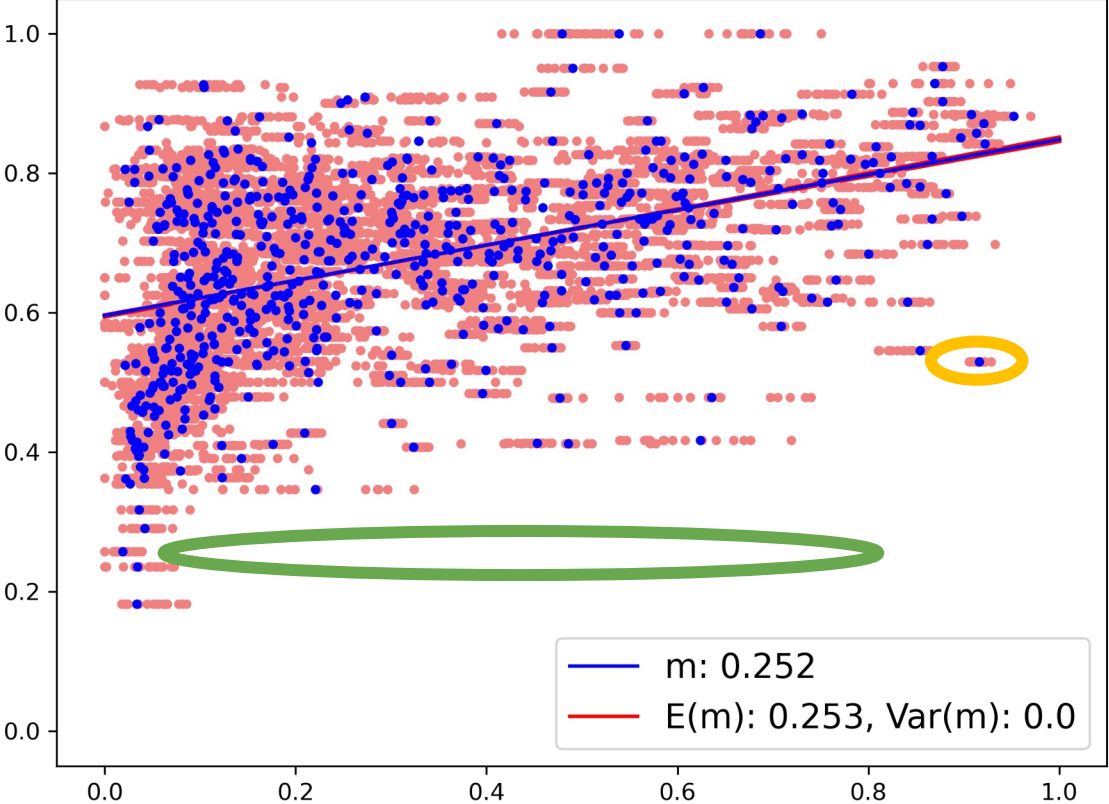
**Pink = Noised Numbers**

# Accounting for small precincts

### Standard ecological regression



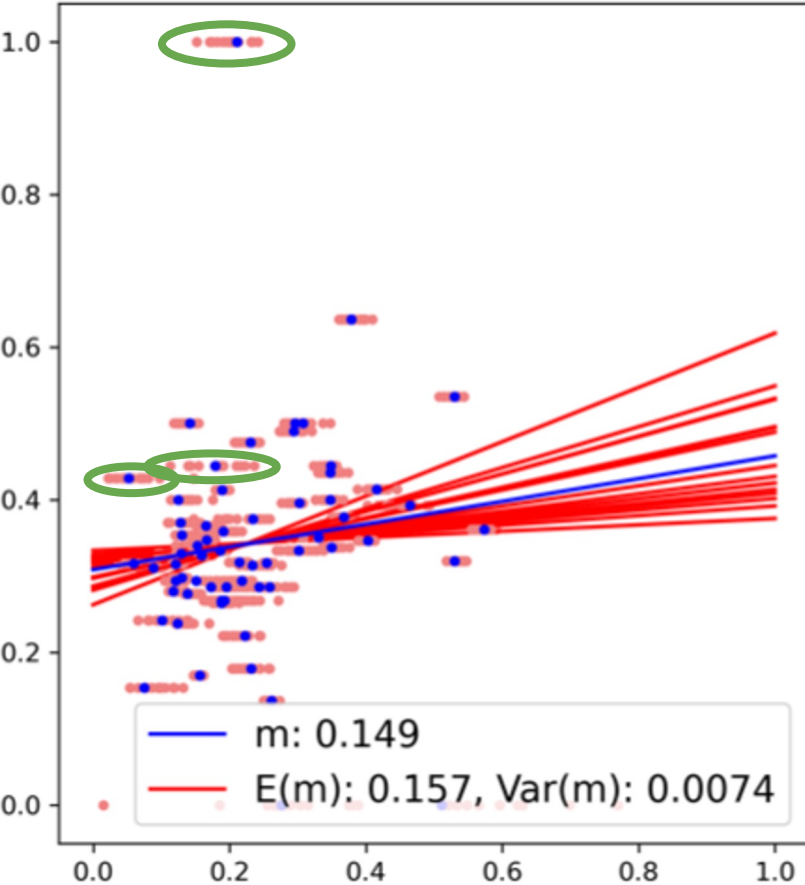
### Filtering precincts with < 10 votes (24% of precincts!)



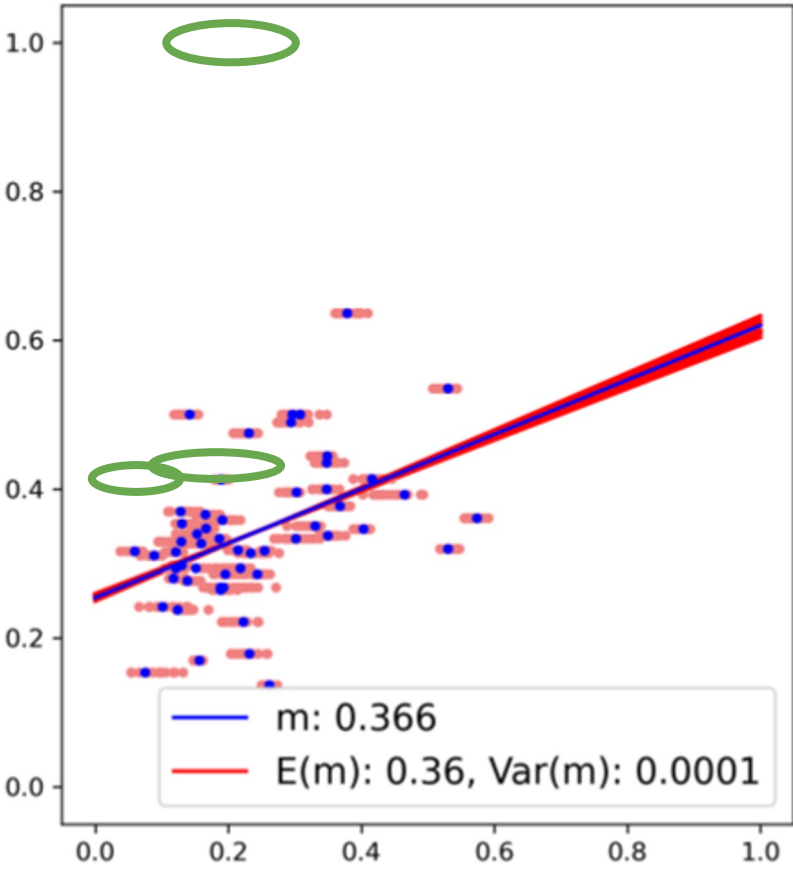
% Hispanic Population  
(Census data)

# Accounting for small precincts

Standard ecological regression



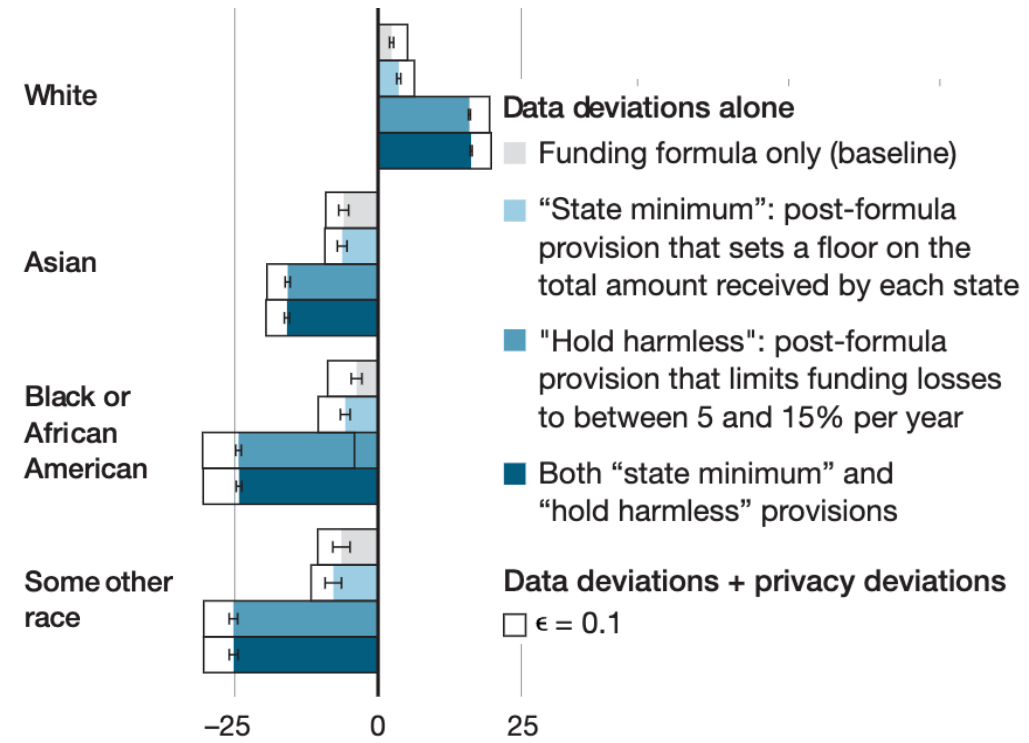
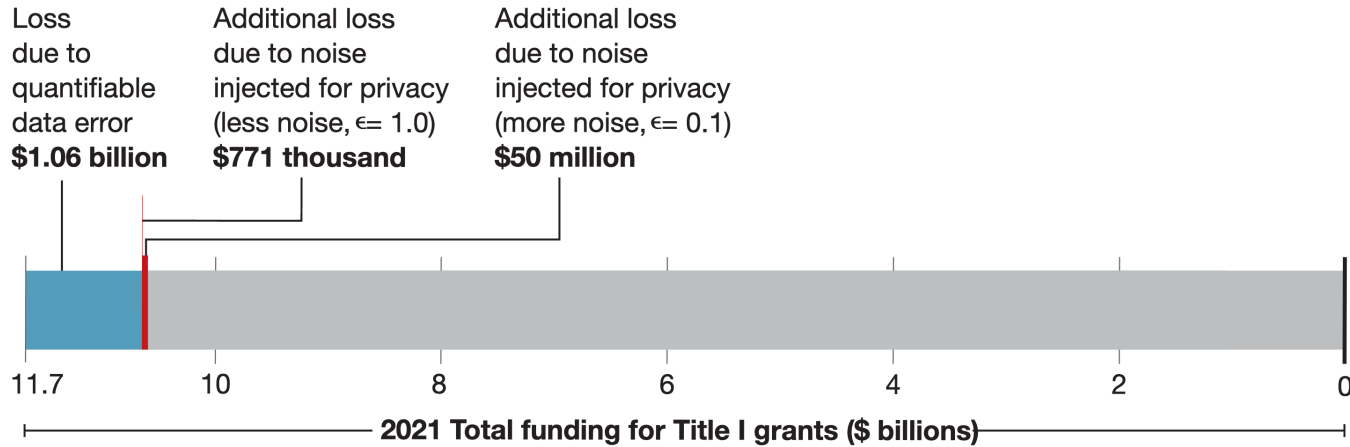
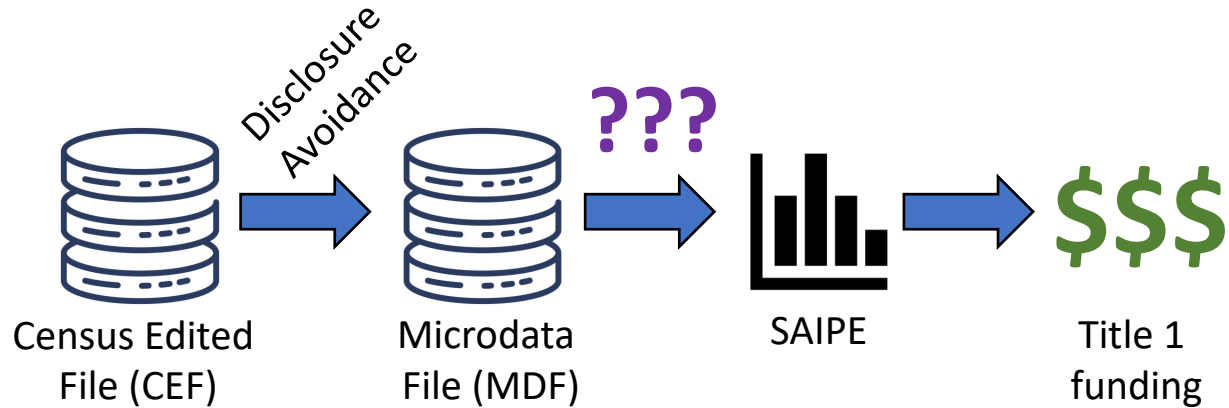
Filtering precincts with < 10 votes



% Support for Valdez  
(Election data)

% Hispanic Population  
(Census data)

# Option 4: Go straight to the policy question



Race-weighted misallocation per eligible child (\$) [edited, partial figure]



**UNITED STATES DISTRICT COURT FOR THE  
MIDDLE DISTRICT OF ALABAMA  
EASTERN DIVISION**

10 10 P 12 37  
THE STATE OF ALABAMA; ROBERT  
ADERHOLT, Representative for Alabama's  
4th Congressional District, in his official and  
individual capacities; WILLIAM GREEN;  
and CAMARAN WILLIAMS,

Plaintiffs,

v.

UNITED STATES DEPARTMENT OF  
COMMERCE; GINA RAIMONDO, in her  
official capacity as Secretary of Commerce;  
UNITED STATES BUREAU OF THE  
CENSUS, an agency within the United States  
Department of Commerce; and RON  
JARMIN, in his official capacity as Acting  
Director of the U.S. Census Bureau,

Defendants.

CIVIL ACTION NO. \_\_\_\_\_

COMPLAINT FOR DECLARATORY AND  
INJUNCTIVE RELIEF

**THREE-JUDGE COURT REQUESTED  
PURSUANT TO 28 U.S.C. § 2284**

## INTRODUCTION

1. This suit challenges two unlawful actions by the U.S. Commerce Department and Census Bureau in relation to the 2020 decennial census—(1) Defendants’ decision to produce manipulated redistricting data to the States, and (2) Defendants’ refusal to produce redistricting data on time.

2. First, the skewed numbers. Congress has ordered the Secretary of Commerce to work with the States to learn what they need for redistricting and then report to each State accurate “[t]abulations of population” for subparts of each State for purposes of “legislative apportionment or districting of such State.” 13 U.S.C. § 141(c). But the Secretary, through the Census Bureau, has announced that she will instead provide the States purposefully flawed population tabulations. The Bureau intends to use a novel statistical method called differential privacy to intentionally skew the population tabulations provided to States to use for redistricting. Thus, while the Bureau touts its mission “to count everyone once, only once, and in the right place,”<sup>1</sup> it will force Alabama to redistrict using results that purposefully count people in the wrong place.

BRIEF OF *AMICI CURIAE* STATE OF  
UTAH AND 15 OTHER STATES IN  
SUPPORT OF PLAINTIFFS

The States of Utah, Alaska, Arkansas, Florida, Kentucky, Louisiana, Maine, Mississippi, Montana, Nebraska, New Mexico, Ohio, Oklahoma, South Carolina, Texas, and West Virginia (*Amici* States) agree with Plaintiffs that the Secretary's intended use of differential privacy deprives states of accurate "[t]abulations of population" of state subparts to use in legislative apportionment and districting under 13 U.S.C. § 141(c). *Amici* States also agree that the Secretary can

# Legal question 1: Was swapping private enough?

- **Law**

- The Census Bureau may not “disclose the information reported by, or on behalf of, any particular respondent”

- **Alabama**

- “there is not a single documented case of anyone outside the Census Bureau revealing the responses of a particular identified person in public use decennial census”

- **Census Bureau**

- “swapping and top and bottom coding applied at the level used in the 2010 census are insufficient to prevent re-identification given the ability to perform database reconstruction and the availability of external data.”

## Legal question 2:

Does TopDown violate “one person, one vote”?

- **Law**

- Congressional districts “as nearly of equal population as is practicable”  
- *Reynolds v Sims (1964)*

- **Alabama**

- "Congressional districts drawn from the demonstration data would likely violate one-person, one-vote"

- **Census**

- "[T]he ‘good-faith effort to achieve population equality’ required of a State conducting intrastate redistricting does not translate into a requirement that the Federal Government conduct a census that is as accurate as possible." - *Wisconsin v New York (1996)*



# Legal question 3: Is TopDown an illegal “statistical method”?

- **Law**

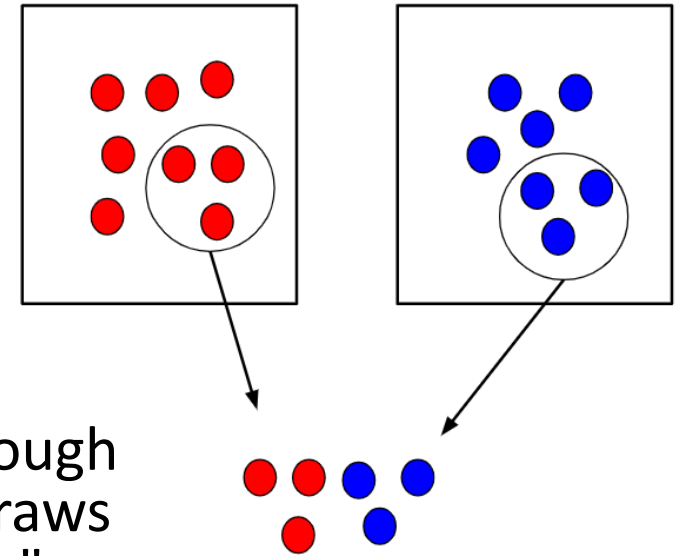
- "the term ‘statistical method’ means ... [any] statistical procedure ... to add or subtract counts to or from the enumeration of the population as a result of statistical inference"

- **Alabama**

- "Privacy is introduced ... by introducing random error through sampling from statistical distributions.... These random draws are then added or subtracted to the actual observations..."

- **Census**

- DP not "a result of statistical inference", and "thus unlike the sampling" at issue in precedent



IN THE UNITED STATES DISTRICT COURT  
FOR THE MIDDLE DISTRICT OF ALABAMA  
EASTERN DIVISION

THE STATE OF ALABAMA, *et al.*,

Plaintiffs,

v.

UNITED STATES DEPARTMENT  
OF COMMERCE, *et al.*,

Defendants.

)  
)  
)  
)  
) Case No. 3:21-cv-211-RAH-ECM-KCN  
) (WO)  
)  
)  
)  
)

**MEMORANDUM OPINION AND ORDER**

After the benefit of oral argument, the court concludes that Plaintiffs' motion for a preliminary injunction and petition for writ of mandamus are due to be DENIED.

Translation:  
Census wins

# The outcome had little to do with DP

Alabama needs **standing** to sue

- Injury in fact
- Traceable
- Redress
- Right of Action

Do individual voters have standing?

- Some voters' power diluted, some amplified
- Catch 22: Can't know which is which.





# Future directions

- More work (empirical + theory) needed on...
  - Reconstruction & privacy
  - Downstream policy impacts
  - "Noisy measurement files"
- What's next from Census?  
→ Next talk!

