

# **Modernizing Disclosure Avoidance: Report on the 2020 Disclosure Avoidance Subsystem as Implemented for the 2018 End-to-End Test (Continued)**

Simson L. Garfinkel  
Chief, Center for Disclosure Avoidance Research  
U.S. Census Bureau

2017 Census Scientific Advisory Committee Fall Meeting  
Suitland, MD  
11:00AM  
September 15, 2017

# Acknowledgments

---

This presentation incorporates work by:

- Dan Kifer (Scientific Lead)
- John Abowd (Chief Scientist)
- Tammy Adams, Robert Ashmead, Aref Dajani, Jason Devine, Michael Hay, Cynthia Hollingsworth, Meriton Ibrahimi, Michael Ikeda, Philip Leclerc, Ashwin Machanavajjhala, Christian Martindale, Gerome Miklau, Brett Moran, Ned Porter, Anne Ross and William Sexton

# Outline

---

Motivation

Differentially private 2020 Disclosure Avoidance System

High-level goals

Flow diagrams

Query examples

Conclusion

# Motivation:

## To protect the privacy of individual survey responses

---

### 2010 Census:

- 5.6 billion independent tabular summaries published.
- Based on 308 million person records

Database reconstruction (Dinur and Nissim 2003) is a serious disclosure threat that all statistical tabulation systems from confidential data must acknowledge.

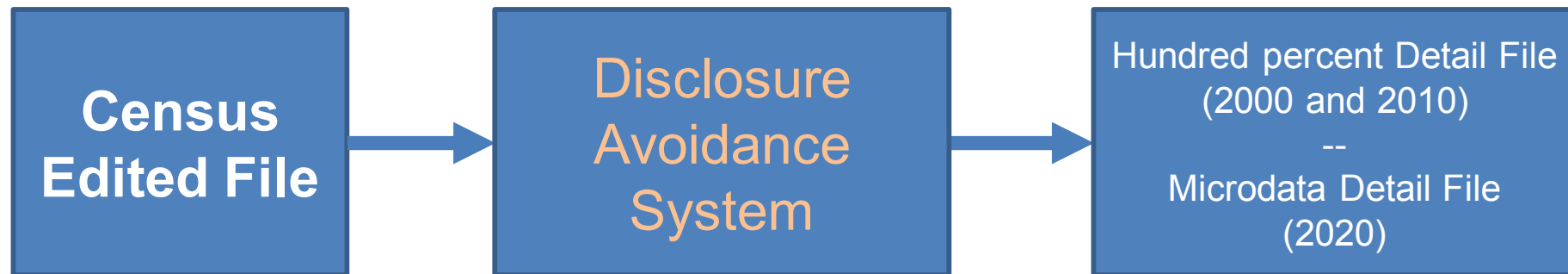
The confidentiality edits applied to the 2010 Census were not designed to defend against this kind of attack.

# The Disclosure Avoidance Subsystem (DAS) implements the privacy protections for the decennial Census.

---

## Features of the DAS:

- Operates on the edited Census records
- Designed to make Census records safe to tabulate



# The 2000 and 2010 Disclosure Avoidance Systems relied on swapping households:

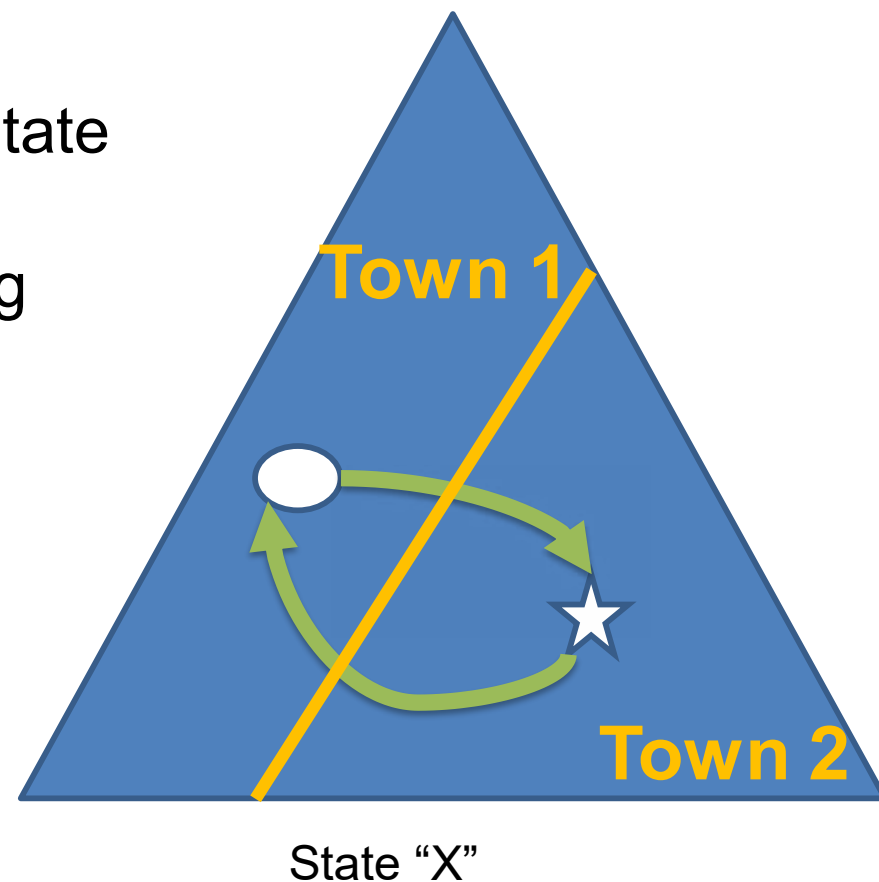
---

## Advantages of swapping:

- Easy to understand
- Does not affect state counts if swaps are within a state
- Can be run state-by-state
- Operation is “invisible” to rest of Census processing

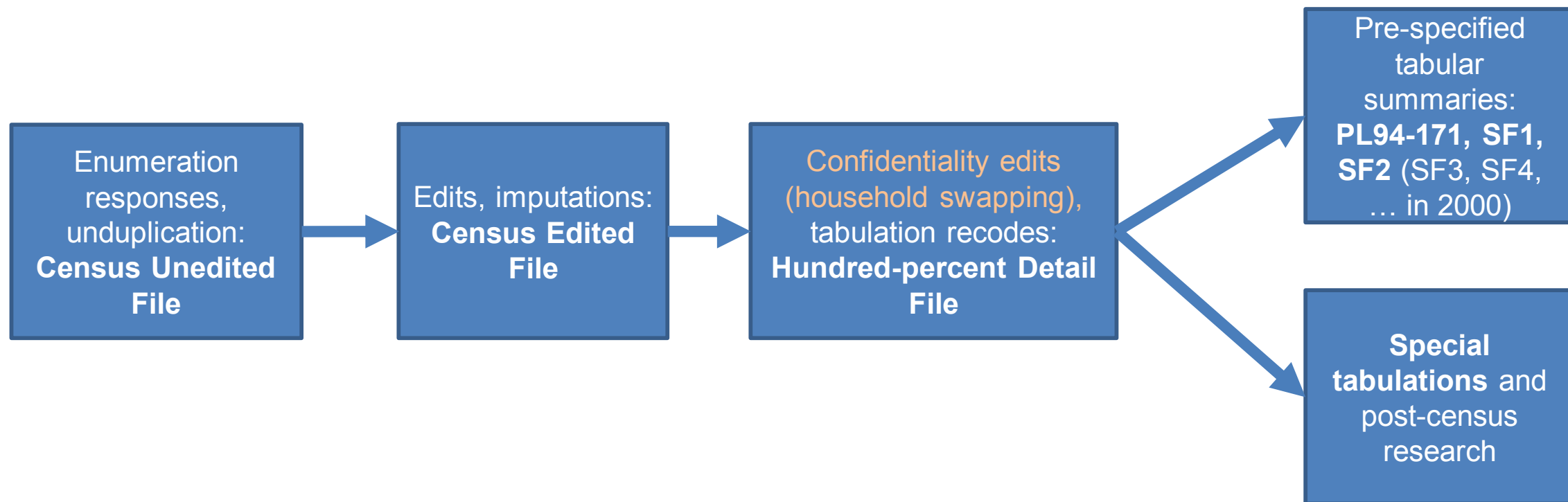
## Disadvantages:

- Does not provide formal privacy guarantees
- Does not protect against database reconstruction attacks
- Privacy guarantee relies on lack of external data



# The 2000 and 2010 Disclosure Avoidance System operated as a filter, on the Census Edited File:

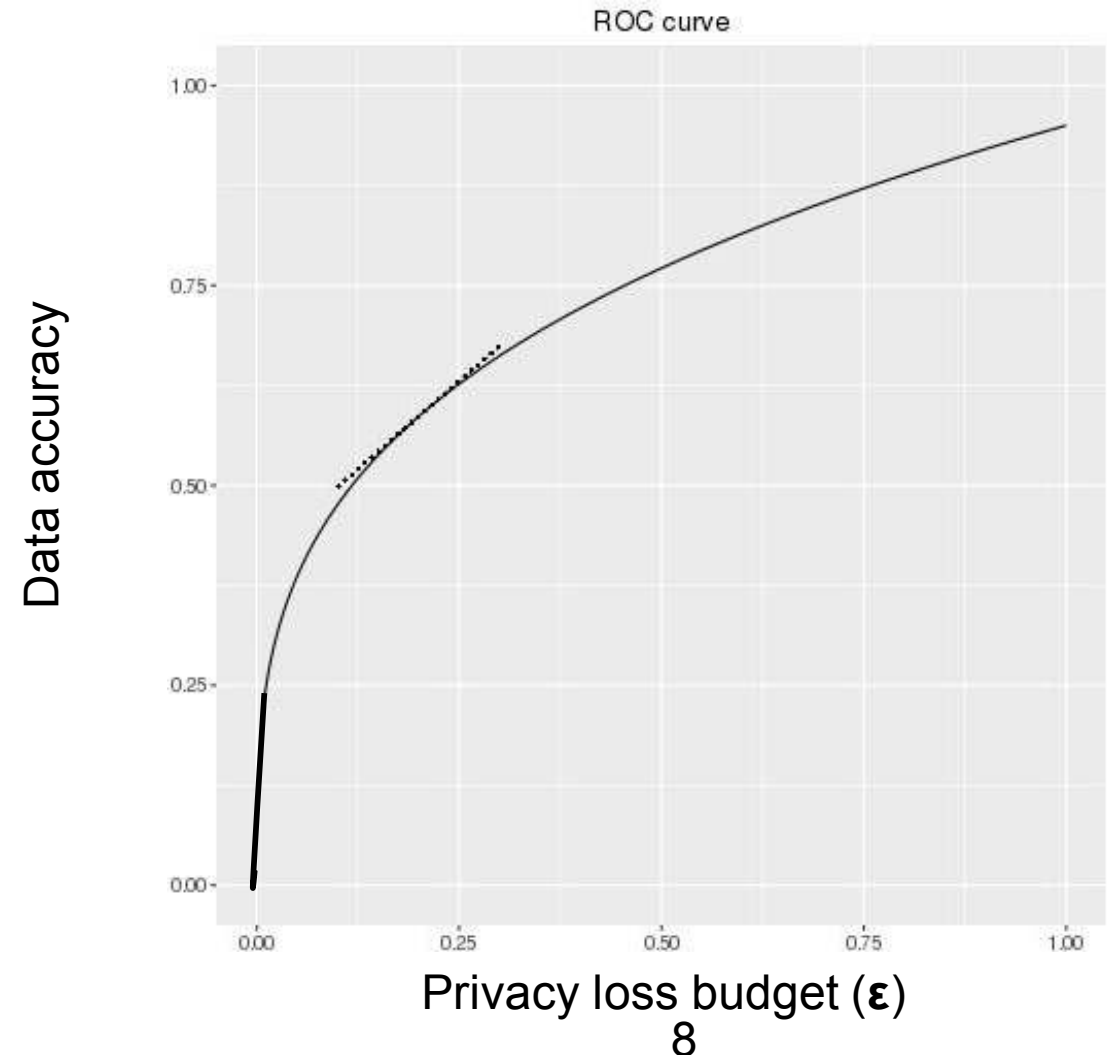
---



# The 2020 Census disclosure avoidance system will use differential privacy to defend against a reconstruction attack,

Differential privacy provides:

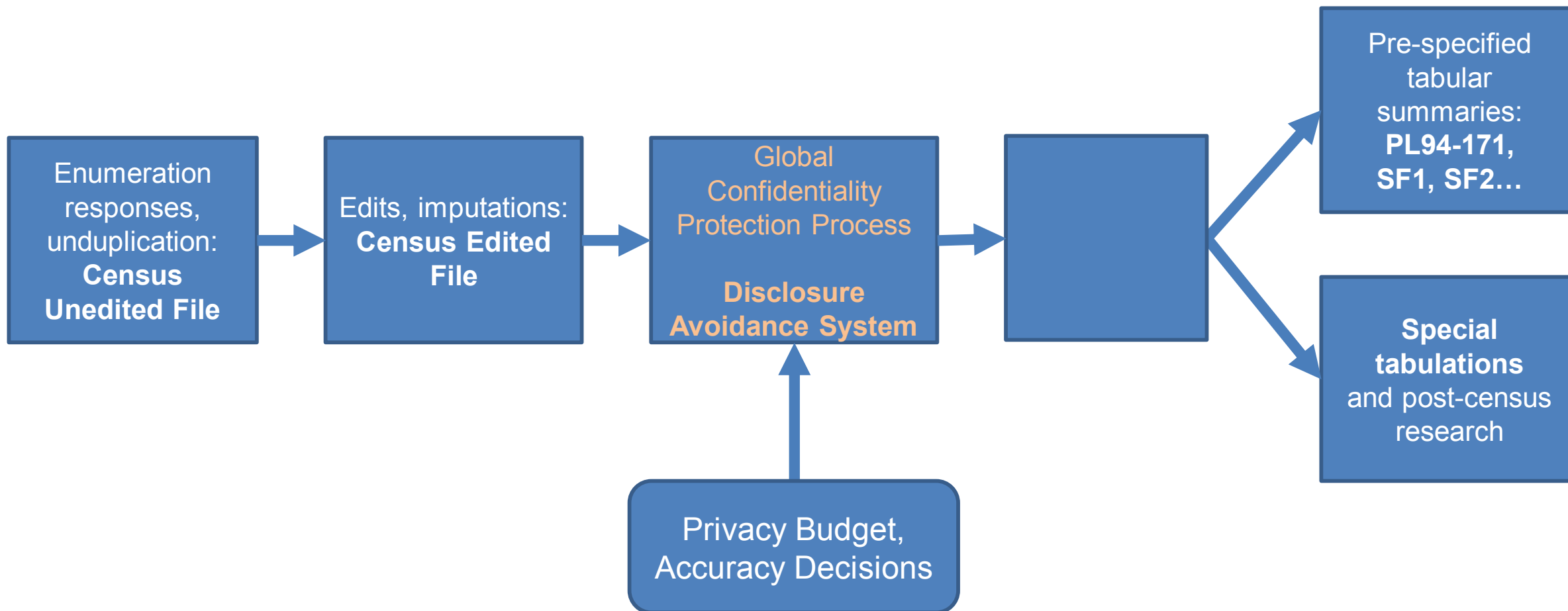
- Provable bounds on the accuracy of the best possible database reconstruction given the released tabulations.
- Algorithms that allow policy makers to decide the trade-off between accuracy and privacy.





# The 2020 Disclosure Avoidance System also operates as a filter... but it's much more visible.

---

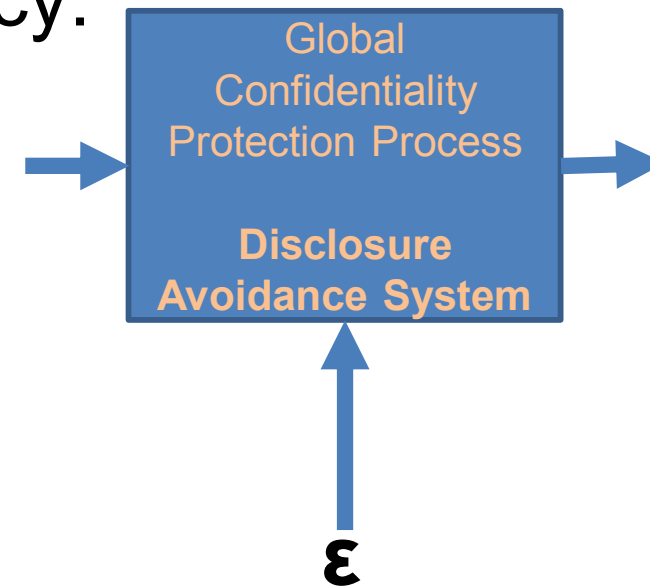


# The 2020 Disclosure Avoidance System relies on infusing formally private noise.

---

Advantages of noise infusion with formal privacy:

- Easy to understand
- Provable and *tunable* privacy guarantees
- Privacy guarantees do not depend on external data
- Protects against database reconstruction attacks
- Privacy operations are *composable*



Disadvantages:

- Entire country must be processed at once for best accuracy
- Every use of private data must be tallied in the *privacy loss budget*

# Differentially Private Disclosure Avoidance System: Requirements

---

DAS must be able to read the Census Edited File (CEF):

- CEF must be exactly specified and contain all information necessary for all tabulation recodes
- CEF must be kept confidential after DAS runs (as it was for historical censuses)

DAS must generate the Microdata Detail File (MDF):

- Must contain all information that appears in *any publicly released table* (e.g. PL94-171, SF1, SF2)
- Should not contain *any information* that does *not* appear in a publicly released table
- May be publicly released (in whole or in part)

Non-functional requirements:

- The disclosure avoidance system must provably move information from the CEF to PL94/SF1/SF2 with an adjustable total privacy-loss budget
- The source code and parameters for the DAS will be made publicly available

# Why generate a differentially private MDF?

---

- Familiar to internal and external stakeholders
- Operates with tabulation system to produce PL-94 and SF-1 tabulations
- Guarantees population totals (voting age, non-voting age, householder) exact at all levels of geography
- Consistency among query answers

# Some queries must be privacy preserving. Some queries must be exact (“invariant”)

Specific PL-94 queries must be exact:

- Block population
- Block voting age population
- Block householders & vacancies

*—per 2000 Department of Justice letter to the Director, Kenneth Prewitt*

Other PL-94 and SF-1 queries will not be exact:

- Age distribution under 18
- Age distribution 18 and over
- Race and ethnicity distribution
- Household relationship distribution
- Household ownership distribution

Final privacy-loss budget determined by Data Stewardship Executive Policy Committee (DSEP) with recommendation from Disclosure Review Board (DRB)

# How the 2020 System Works:

## High-level Overview

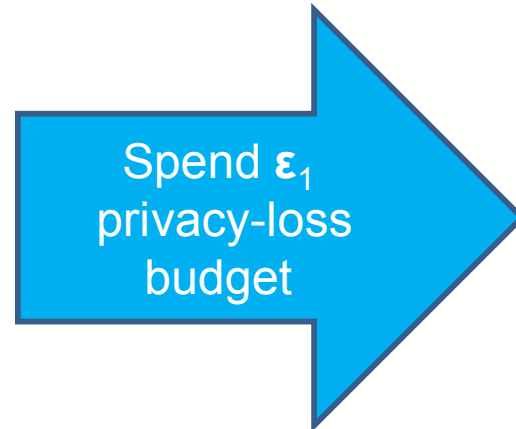
---

The new system is similar to swapping, with key differences:

- Every record in the population may be modified  
*But modifications are bounded by DOJ policy and global privacy budget.*
- Records in the tabulation data have no exact counterpart in the confidential data  
*There is no one-to-one mapping between CEF and MDF records.  
But there are the same number of records for every block.*
- Explicitly protected tabulations (PL-94 and SF-1) have provable, public accuracy levels  
*2020 will publish the algorithms, the parameters and the accuracy of the tabulations.*

# Proposed “Top-Down” Algorithm

National table of  
US population  
 $2 \times 255 \times 17 \times 115$



National table with all 500,000 cells  
filled, structural zeros imposed with  
accuracy allowed by  $\epsilon_1$   
 $2 \times 255 \times 17 \times 115$



Reconstruct individual micro-data  
without geography  
325,000,000 records



Sex: Male / Female  
Race + Hispanic: 255 possible values  
Relationship to Householder: 17  
Age: 0-114



# State-level

State-level tables for only certain queries; structural zeros imposed; dimensions chosen to produce best accuracy for PL-94 and SF-1

Spend  $\epsilon_2$   
privacy-loss  
budget

Target state-level tables required for best accuracy for PL-94 and SF-1  
Exact state voting-age, non-voting age, and householder counts as enumerated.



Construct best-fitting individual micro-data with state geography

325,000,000 records now including state identifiers



# County-level

County-level tables for only certain queries; structural zeros imposed; dimensions chosen to produce best accuracy for PL-94 and SF-1

Spend  $\epsilon_3$   
privacy-loss  
budget

Target county-level tables required for best accuracy for PL-94 and SF-1

Exact county voting-age, non-voting age, and householder counts as enumerated.



Construct best-fitting individual micro-data with state and county geography

325,000,000 records now including state and county identifiers

# Census tract-level

Tract-level tables for only certain queries; structural zeros imposed; dimensions chosen to produce best accuracy for PL-94 and SF-1

Spend  $\epsilon_4$   
privacy-loss  
budget

Target **tract-level** tables required for best accuracy for PL-94 and SF-1

Exact **tract** voting-age, non-voting age, and householder counts as enumerated.



Construct best-fitting individual micro-data with **state, county, and tract** geography

325,000,000 records now including **state, county, and tract** identifiers

# Block-level

Block-level tables for only certain queries; structural zeros imposed; dimensions chosen to produce best accuracy for PL-94 and SF-1

Spend  $\epsilon_5$  privacy-loss budget

Block tract-level tables required for best accuracy for PL-94 and SF-1  
Exact block voting-age, non-voting age, and householder counts as enumerated.



Construct best-fitting individual micro-data with state, county, tract and block geography  
325,000,000 records now including state, county, tract identifiers

# MDF for tabulating

Construct best-fitting individual micro-data  
with **state, county, tract and block**  
geography

325,000,000 records now including state,  
county, tract, and block identifiers



MDF used for tabulating  
PL-94, SF-1

# MDF for tabulating

How accurate is the MDF?



Disclosure Avoidance Certificate

- Certifies that the DAS passed tests
- Reports the accuracy of the MDF
- Requires  $\epsilon_A$

Construct best-fitting individual micro-data  
with **state, county, tract and block**  
geography

325,000,000 records now including state,  
county, tract, and block identifiers

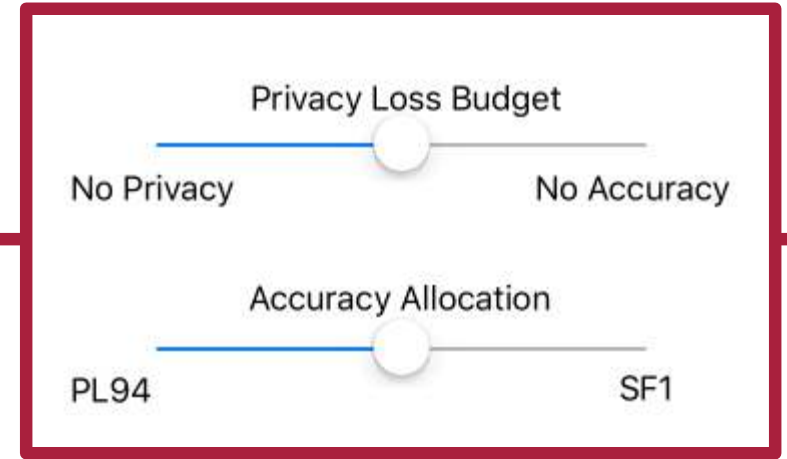


MDF used for tabulating  
PL-94, SF-1

# Operational Decisions

Set total privacy loss budget:  $\epsilon$

- Ensure that  $\epsilon_1 + \epsilon_2 + \epsilon_3 + \epsilon_4 + \epsilon_5 + \epsilon_A = \epsilon$



Within each stage, allocate privacy-loss budget between:

- PL-94
- Parts of SF-1 not in PL-94

These are policy levers provided by the system.

Levers are set by the Data Stewardship Executive Policy Committee

# Inputs Used by the Development Team

---

Lists of matrices in technical documentation express core queries in the workload

- PL94: <https://www.census.gov/prod/cen2010/doc/pl94-171.pdf>
- SF1: <https://www.census.gov/prod/cen2010/doc/sf1.pdf>
- SF2: <https://www.census.gov/prod/cen2010/doc/sf2.pdf>

Over 1,000 pages of edit specifications for 2010 CEF

Uncurated tabulation recode programs

# We are creating

## A framework for Disclosure Avoidance Systems:

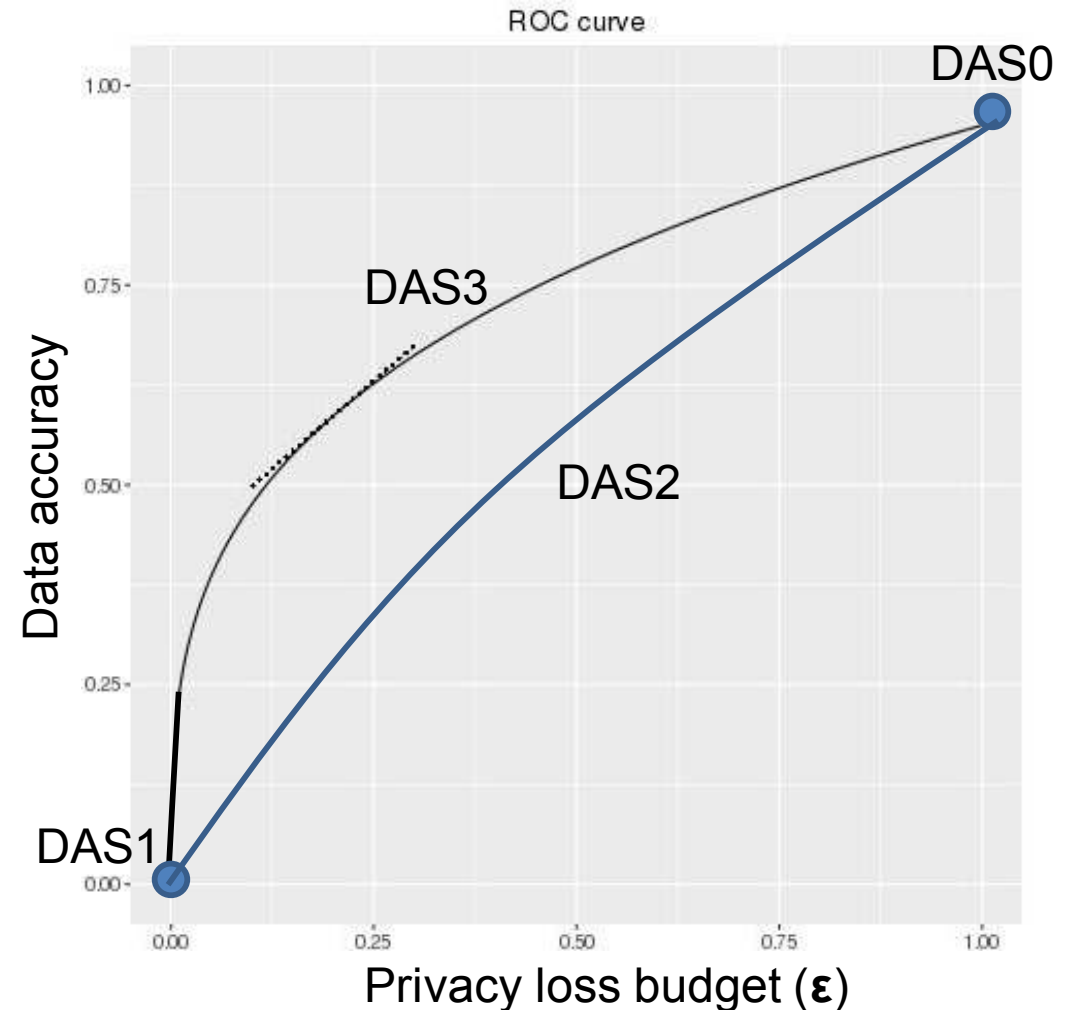
- Development & Test Mode
- Production Mode

## Testing Systems:

- DAS0 — 100% accuracy, no privacy  
(No disclosure avoidance)
- DAS1 — 100% privacy, no accuracy
- DAS2 — “bottom-up” engine

## Operational System:

- DAS3 — “top-down” engine





# Plans for the 2018 End-to-End Test

---

The 2018 End-to-End test will incorporate differential privacy

- Likely DAS2 — Bottom-up algorithm

Only the prototype PL94-171 files will be produced

No decisions yet regarding the privacy-loss budget or accuracy level

Questions?

# Reference

---

Dinur, Irit and Kobbi Nissim (2003). “Revealing information while preserving privacy.” in *Proceedings of the twenty-second ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems* (PODS '03). ACM, New York, NY, USA, 202-210. DOI: 10.1145/773153.773173.