# Updating Data Editing and Review for Large-scale Commodity Flow Survey Response Data

2023 Federal Committee on Statistical Methodology
October 24, 2023
College Park Marriott Hotel & Conference Center
3501 University Blvd E, Hyattsville, MD 20783

Gritiya Tanner

United States® Census Bureau

# Disclaimers

Any conclusions expressed herein are those of the authors and do not reflect the views of the U.S. Census Bureau.

The Census Bureau has reviewed this data product to ensure appropriate access, use, and disclosure avoidance protection of the confidential source data (Project No. P-7504831, Disclosure Review Board (DRB) approval number:  CBDRB-FY24-ESMD002-002).

# Background

Challenges:

- 16 times data increase
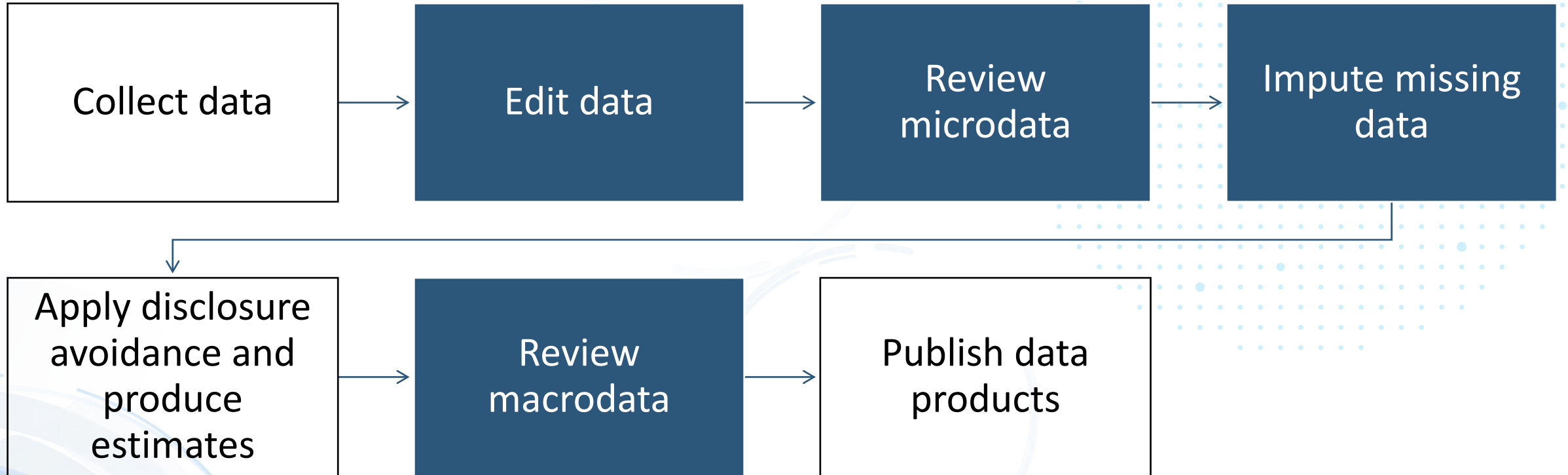- Computational costs
- Storage requirements
- Data review

Solution:
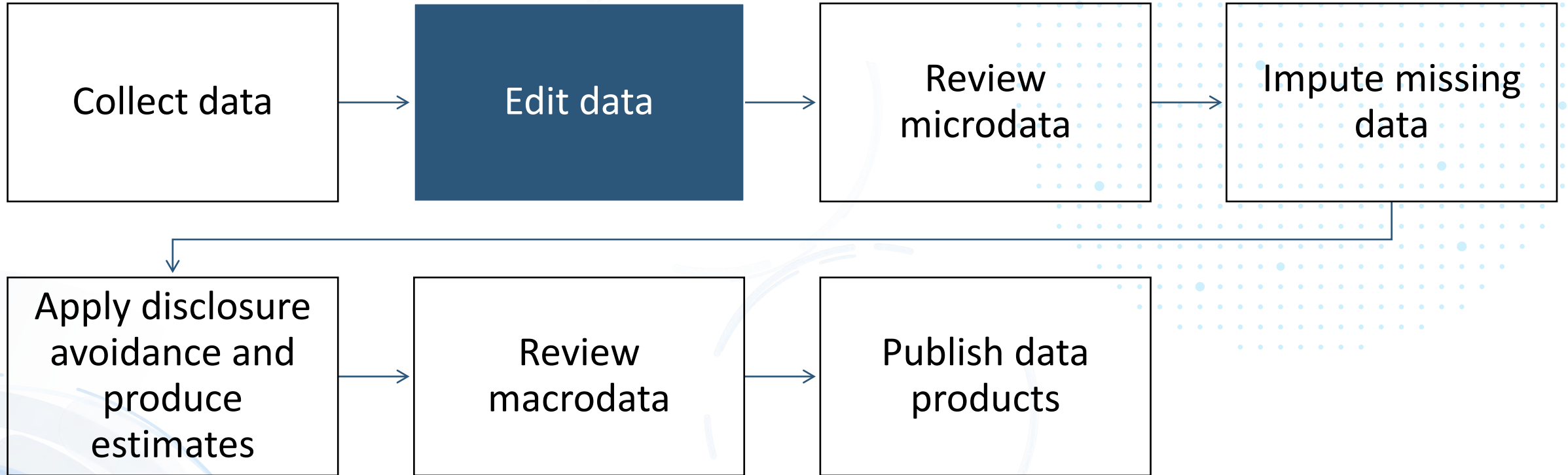
- Redesign and automate the process as much as you can

Today:

- Highlight the implemented changes from 2017 CFS to 2022 CFS

# 2022 CFS Data processing pipeline

Collect data → Edit data → Review microdata → Impute missing data

Apply disclosure avoidance and produce estimates → Review macrodata → Publish data products

United States® Census Bureau

4

# 2022 CFS Data processing pipeline

```
Collect data → Edit data → Review microdata → Impute missing data
                                                        ↓
Apply disclosure avoidance and produce estimates → Review macrodata → Publish data products
```

# Data editing process

**Edit flag to identify:**

- Invalid data
- Missing data
- Common data issues

**Type of edit flag:**

- Shipment
- Establishment
- Quarter-to-quarter
- Hazmat
- Imputation

**Data correction:**

- Manually review each record and make changes through the user interface
- Programmatically apply changes based on set parameters through batch collection tool

# Changes in data editing process
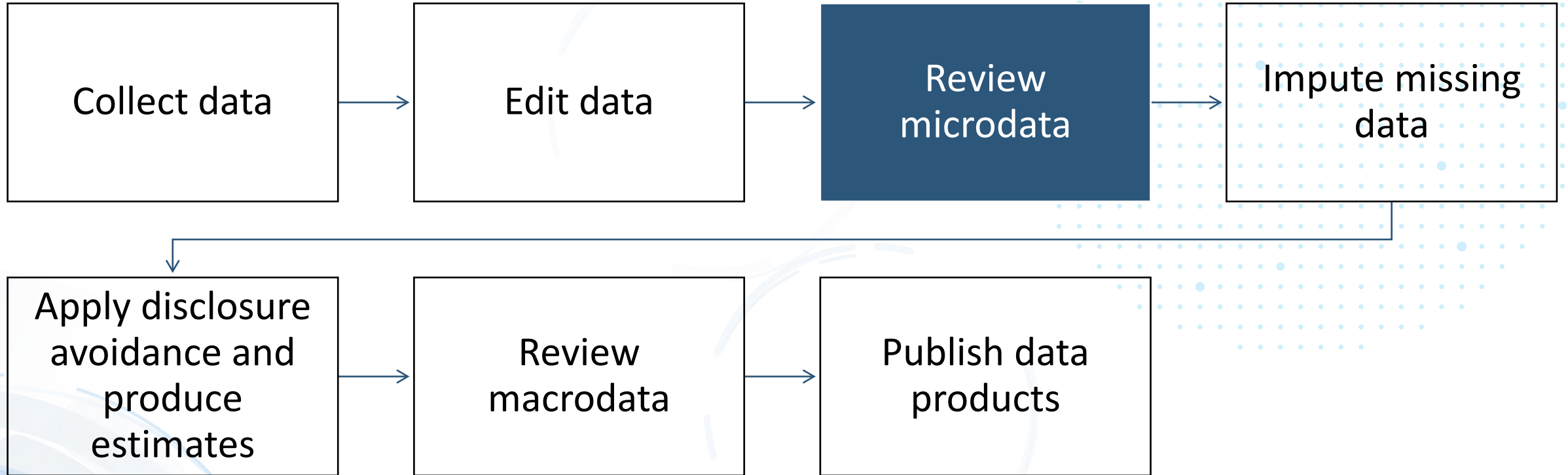
## 2017 CFS

- **Data**
  - 6.5 million shipments
  - 7 million edit failures
- **Edit checks**
  - 3 categories
  - Run every weekend
- **Batch data correction**
  - Ad-hoc process
  - Take 1 week to reflect data changes in database
- **Data review tools**
  - Tables

## 2022 CFS

- **Data**
  - 106 million shipments
  - 170 million edit failures
- **Edit checks**
  - 5 categories
  - Run on demand
- **Batch data correction**
  - Streamlined process
  - Take 1 day to reflect data changes in database
- **Data review tools**
  - Data visualization

United States® Census Bureau

# 2022 CFS Data processing pipeline

```
┌─────────────┐     ┌─────────────┐     ┌─────────────┐     ┌─────────────┐
│             │     │             │     │   Review    │     │   Impute    │
│ Collect data│ ──▶ │  Edit data  │ ──▶ │  microdata  │ ──▶ │missing data │
│             │     │             │     │             │     │             │
└─────────────┘     └─────────────┘     └─────────────┘     └─────────────┘

┌─────────────┐     ┌─────────────┐     ┌─────────────┐
│   Apply     │     │             │     │             │
│ disclosure  │     │   Review    │     │Publish data │
│avoidance and│ ──▶ │  macrodata  │ ──▶ │  products   │
│  produce    │     │             │     │             │
│  estimates  │     │             │     │             │
└─────────────┘     └─────────────┘     └─────────────┘
```

# Microdata review process

**Review:**

- Visualization tool

- Prioritize the review
    - Quarter
    - Make it to tabulation
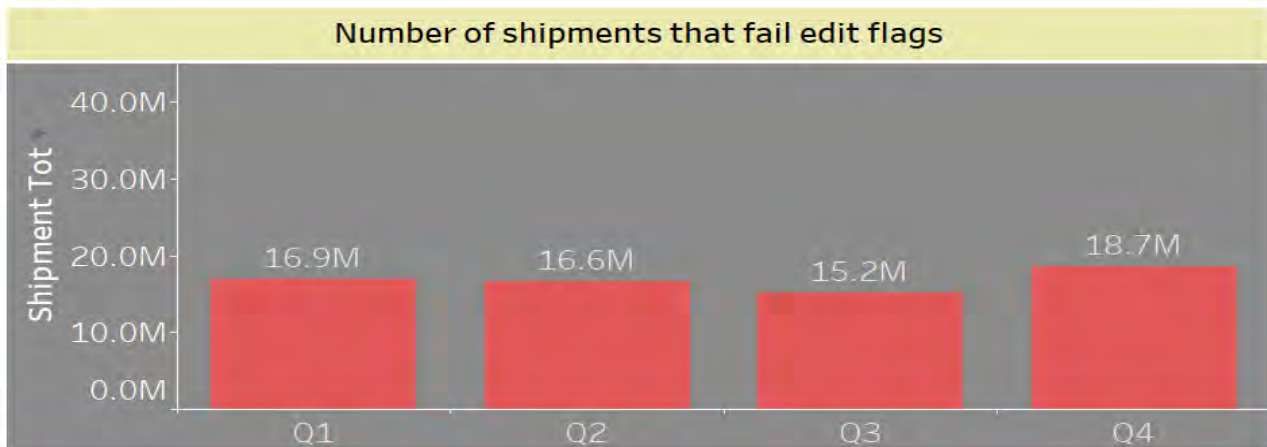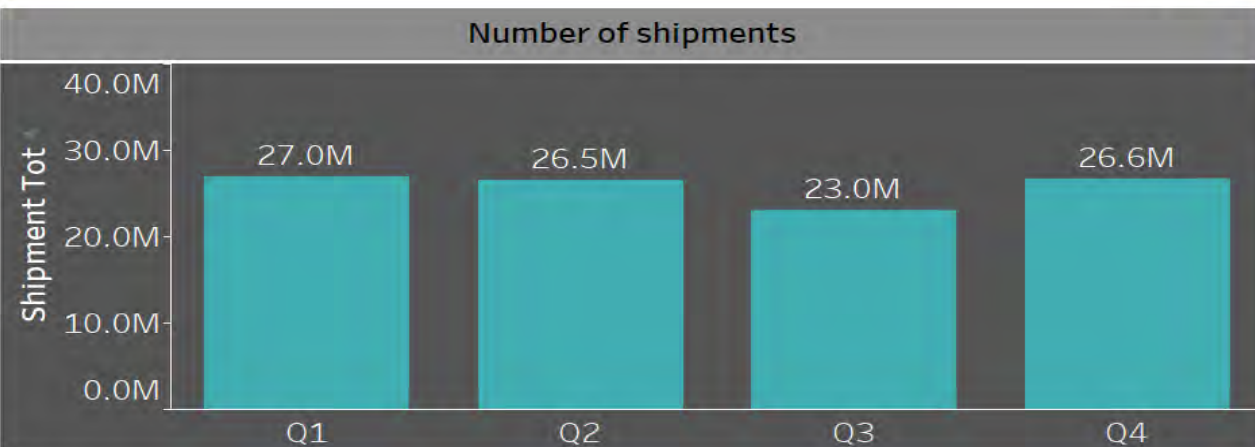
- Figure out the approach to make a correction

**Make any correction or changes:**
- Manually change each record through the user interface
- Programmatically apply changes through batch correction

# Microdata review tool – Visualization summary (New for 2022 CFS)



## 2022 Shipment Flag Metric Dashboard

| Summary | #Shipments by Flag Status | #Shipment by Flag Value | #Flag Per Shipment | Fail 1 flag per shipment | Fail 2 flags or more per shipment |

### Summary

**Number of shipments**

- Q1: 27.0M
- Q2: 26.5M
- Q3: 23.0M
- Q4: 26.6M

**Number of shipments that fail edit flags**

- Q1: 16.9M
- Q2: 16.6M
- Q3: 15.2M
- Q4: 18.7M

**Number of establisments**

- Q1: 44.6K
- Q2: 38.8K
- Q3: 36.0K
- Q4: 35.2K

**Number of fail edit flags**

- Q1: 39.8M
- Q2: 38.8M
- Q3: 35.8M
- Q4: 40.7M

# Microdata review tool – Visualization detail (New for 2022 CFS)

# CFS Production user interface

# Batch correction tool

# 2022 CFS Data processing pipeline

Collect data → Edit data → Review microdata → **Impute missing data**

Apply disclosure avoidance and produce estimates → Review macrodata → Publish data products
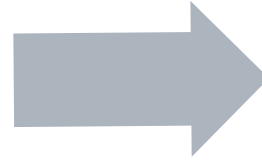
# The changes on imputation process
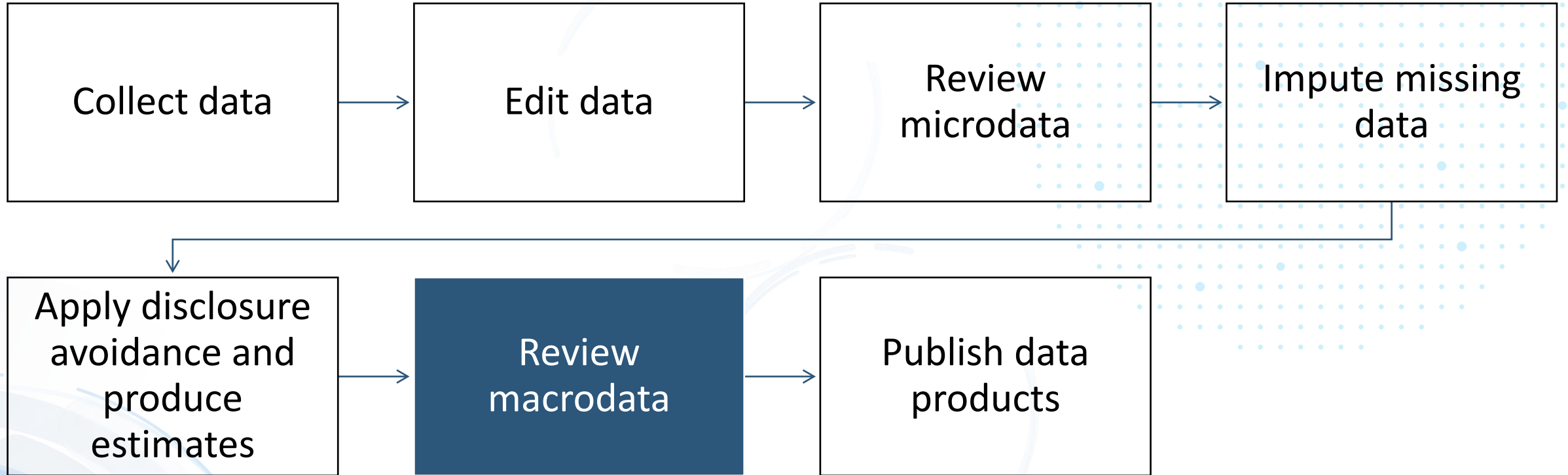
## 2017 CFS

- **Data**
  - 6.5 million shipments
  - 7 million edit failures
- **Imputation process**
  - Only two variables are automatically imputed. The rest are done through ad-hoc imputation.

## 2022 CFS

- **Data**
  - 106 million shipments
  - 170 million edit failures
- **Imputation process**
  - All variables are automatically imputed.
  - Required a database redesign.

# 2022 CFS Data processing pipeline

```
┌──────────────┐      ┌──────────────┐      ┌──────────────┐      ┌──────────────┐
│              │      │              │      │   Review     │      │Impute missing│
│ Collect data │ ───► │  Edit data   │ ───► │  microdata   │ ───► │    data      │
│              │      │              │      │              │      │              │
└──────────────┘      └──────────────┘      └──────────────┘      └──────────────┘

┌──────────────┐      ┌──────────────┐      ┌──────────────┐
│Apply disclosure│    │              │      │              │
│ avoidance and │ ───► │   Review     │ ───► │ Publish data │
│   produce     │      │  macrodata   │      │  products    │
│  estimates    │      │              │      │              │
└──────────────┘      └──────────────┘      └──────────────┘
```

United States®
Census
Bureau

# Macrodata review process

**Review:**

- Visualization tool

- Review and learn macrodata

- Figure out the approach to make a correction

- Perform macrodata review prior to micro data review when estimation available

- Estimation available 1.5 year sooner than last cycle

**Make any correction or changes:**

- Manually review record and make changes through the user interface
- Programmatically apply changes through batch correction

# Macrodata review tool - 2017 CFS version

## EST Compare Dashboard
## HAZ03,Hazardous Materials Series: HazMat Shipment Characteristics by UN/NA Number for the United States: 2017

# Macrodata review tool - new in 2022 CFS



**Shipments Dashboard**

**HAZ03,Hazardous Materials Series: HazMat Shipment Characteristics by UN/NA Number for the United States: 2017**

Back Button

| Value Wt | Value Uwt | Pounds Wt | Pounds Uwt | GCD |
|---|---|---|---|---|
| Tbl Num | Tbl Num | Tbl Num | Tbl Num | Tbl Num |
| HAZ03 | HAZ03 | HAZ03 | HAZ03 | HAZ03 |

Shipment detail

Analyst   Tbl Num   Xtab4 Source   Xtab4   Xtab8 Source   Xtab8   ID   Id Name   Quarter   Line No   MOS   Naics   Sctg   Estab Wt   Shipment Wt

# Summary

To address the challenges associated with the large-scale data processing pipeline, including:

- Automated the process as much as we can
- Redesign the database to support the automated process effectively
- Using the data visualization tool to speed up the data review
- Review macrodata before microdata when possible

These improvements can be adopted by other surveys whether dealing with a large-scale data set or small sets of data.

# Updating Data Editing and Review for Large-scale Commodity Flow Survey Response Data

## QUESTIONS?

## Contacts

**Gritiya Tanner**

Survey Statistician/Data scientist
Business Development Staff

Economic Reimbursable Surveys Division

U.S. Census Bureau

gritiya.tanner@census.gov