# Building AI-Ready Best Practices: The Work of the Commerce Data Governance Board

**October 2024**

*The views expressed in this presentation are those of the presenter and not the Census Bureau.*

**Ryan Harper**, Product Manager

Center for Enterprise Dissemination
U.S. Census Bureau

# DOC AI and Open Government Data Assets Working Group

- The RFI that was sent out on April 17[th], 2024 received over thirty responses

- The working group was established shortly thereafter and features participants from across the DOC

- The initial working group guidance is expected to be released this year

# Generative AI and LLMs

Generative AI applications enable users to create synthetic content, including text, images, audio, and video.

Large language models (LLMs) leverage deep learning algorithms to model language patterns and generate text.

United States® Census Bureau

# Democratizing Access to Public Data Through Generative AI

The emergence of GenAI has the potential to provide improved information and data access to users, from novices to experts.

GenAI and other AI technologies present opportunities and challenges for data providers and data users—including government entities, industry, academia, and the public at large.

United States® **Census** Bureau

# Tech Adoption Patterns Tend to Repeat

Technology emerges

⬇

Public expectations change

⬇

Policy follows

United States® Census Bureau

# Changes in the Public Data Landscape
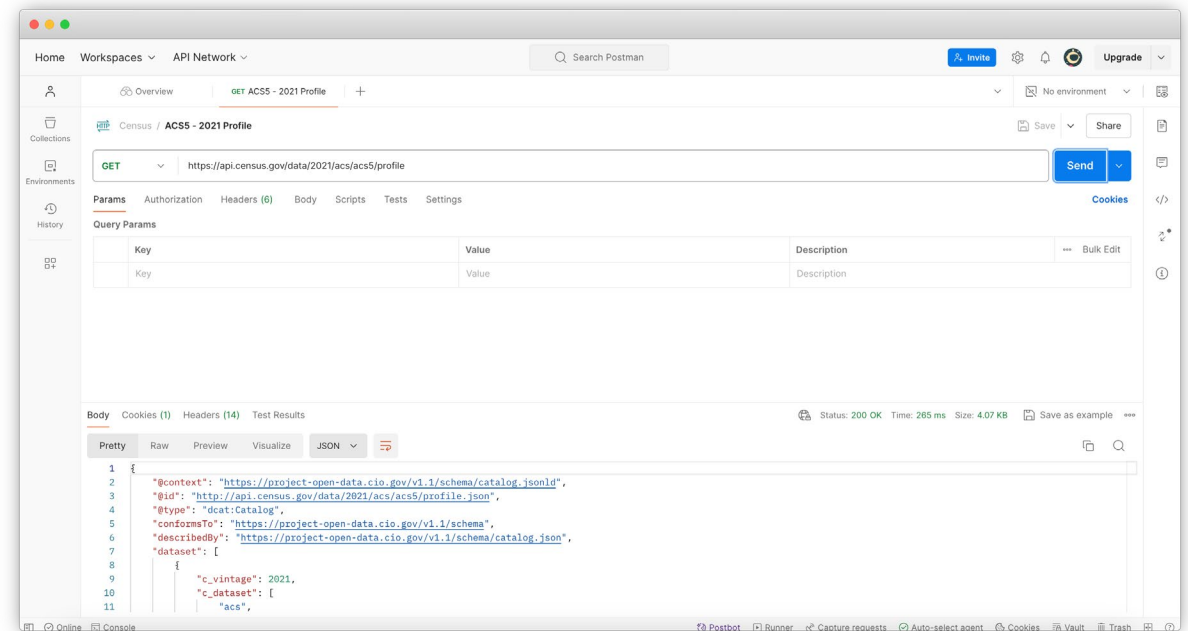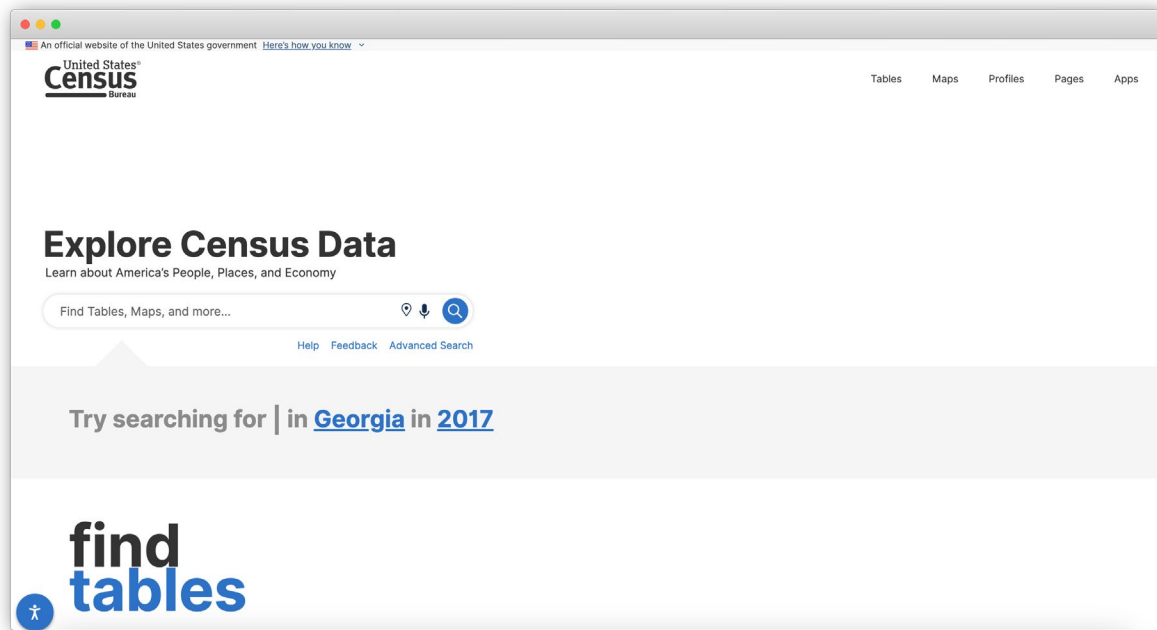
## 20<sup>th</sup> Century

- Federal government was the dominant user

- The statistical system was a near monopoly

- Output was mostly cross-tabs

- Published in books and deposited in libraries, then electronically largely in book formats (PDFs)

- Source data acquisition was difficult and costly

- Privacy and confidentiality risks were small

- Computation was limited and expensive
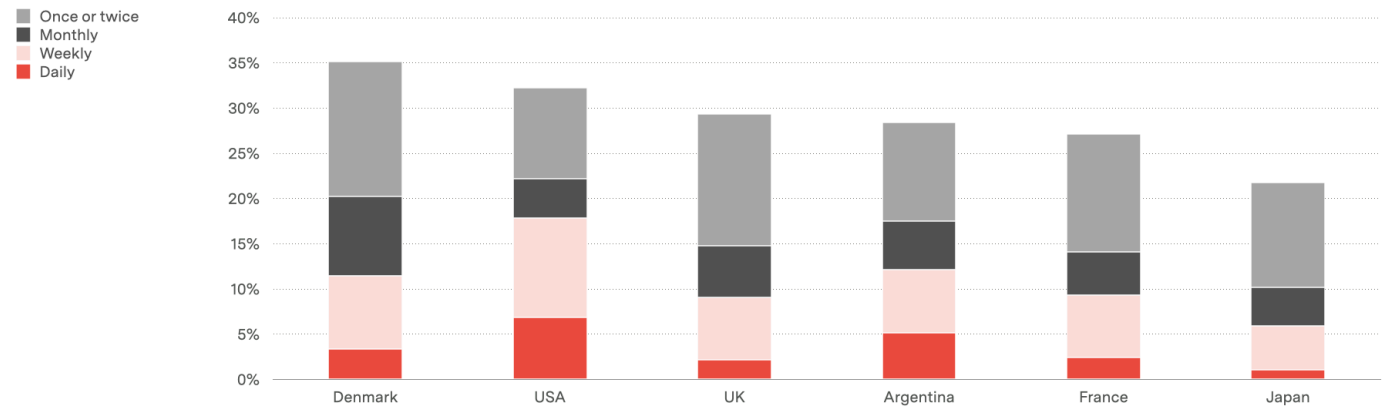
## 21<sup>st</sup> Century

- There are many diverse users of public data

- There are many more organizations that produce similar statistical products

- Output is varied and features visualization and tools for analysis

- Data is accessed electronically online or in secure enclaves

- Source data is more abundant and available in structured and unstructured formats

- Privacy and confidentiality risks are much greater

- Computation is vastly improved

United States® Census Bureau

# data.census.gov and the U.S. Census API

# So Are Generative AI and LLMs Next?

**Use of ChatGPT by country, May 2024**

Legend:
- Once or twice (grey)
- Monthly (dark grey)
- Weekly (light pink)
- Daily (red)

Y-axis: 0% to 40%

Countries: Denmark, USA, UK, Argentina, France, Japan

Source: Reuters Institute

Benedict Evans –– July 2024

United States® Census Bureau

Current State

# ChatGPT, Bard, and CensusGPT return incorrect results from secondary sources (2023)

# ChatGPT, Gemini, and Copilot return incorrect results from secondary sources (February 2024)



**ChatGPT screenshot:**

You
Tell me about the demographics of Suitland, Maryland.

ChatGPT
Suitland is a census-designated place (CDP) in Prince George's County, Maryland, United States, located just outside Washington, D.C. As of the 2020 U.S. Census, Suitland had a population of approximately 26,000 residents. — 25,839 — The racial and ethnic composition of Suitland is diverse, with a significant majority of African American residents, making up around 80% of the population. — 87.4% — White residents account for about 7%, while — 2.5% — Hispanic or Latino residents make up around 10%. — 7.8% — Other racial groups, including Asians and those identifying as two or more races, constitute the remaining percentage of the population.

The median household income in Suitland is lower than the national average, and

**Gemini screenshot:**

Tell me about the demographics of Suitland, Maryland.

Show drafts

Suitland, Maryland, has a population of around 26,400 people as of 2021. — 26,375 — Here's a breakdown of some key demographics:

**Race and Ethnicity:**

- The largest racial/ethnic group is Black or African American (Non-Hispanic), making up around 89% of the population. — 91% — dataus.io
- Hispanic residents account for about 6.7%, — 6.3% — followed by smaller proportions of people identifying as Two or More races, White, and other groups.

Enter a prompt here

Gemini may display inaccurate info, including about people, so double-check its responses. Your privacy & Gemini Apps

gemini.google.com

**Copilot screenshot:**

Demographics of Suitland, ...

You
What are the demographics of Suitland, Maryland?

Copilot

Explore

According to the U.S. Census Bureau, Suitland, Maryland had a population of **25,839** as of April 1, 2020 [1]. The majority of the population was Black or African American (89.5%), followed by Two or More Races (4.8%) and Hispanic or Latino (6.7%) [2]. The median age was **35.8** years, and the median household income was $64,652 [3]. Suitland had a relatively high rate of — 87.4% — 3.8% — 7.8%
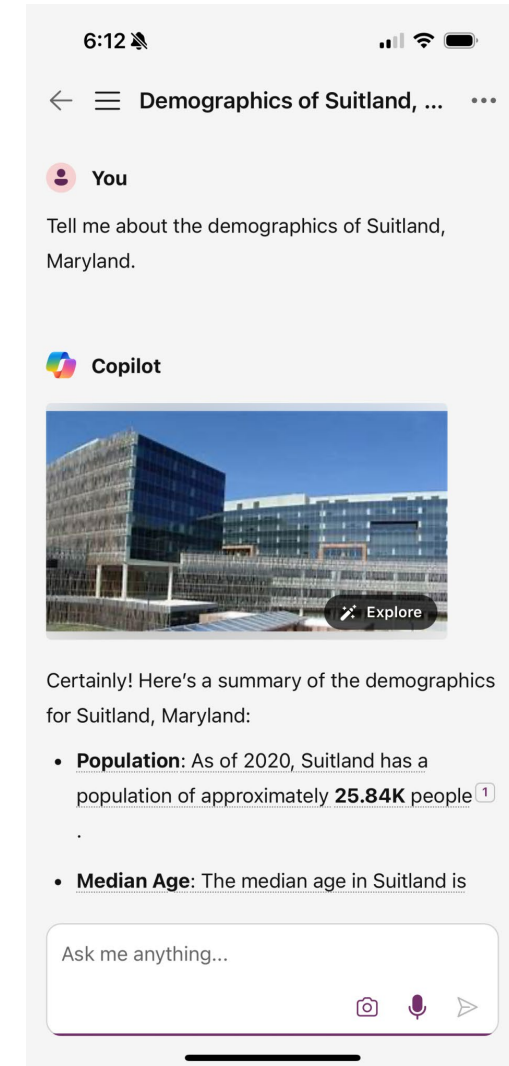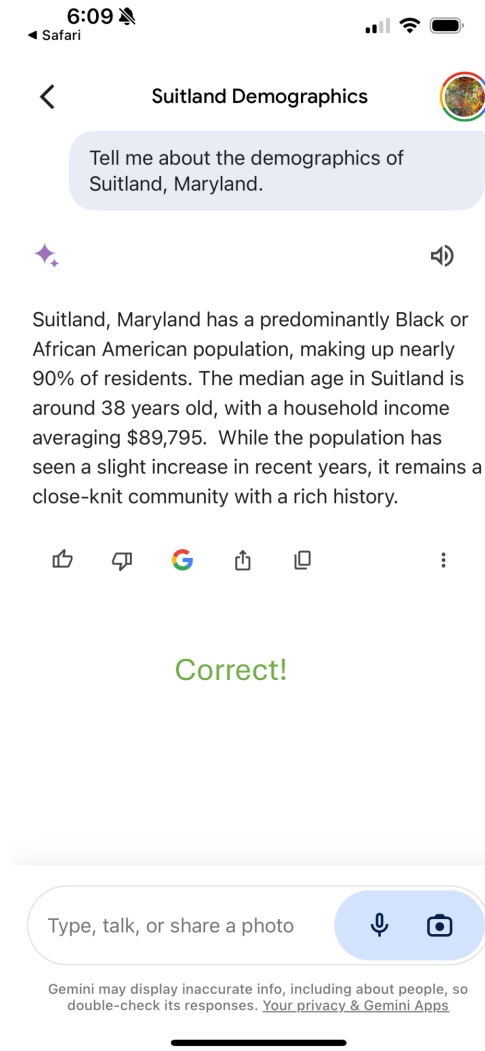
1: neighborhoodscount.com
2: census.gov (QuickFacts)
3: datausa.io

Ask me anything...
0/4000

United States Census Bureau

11
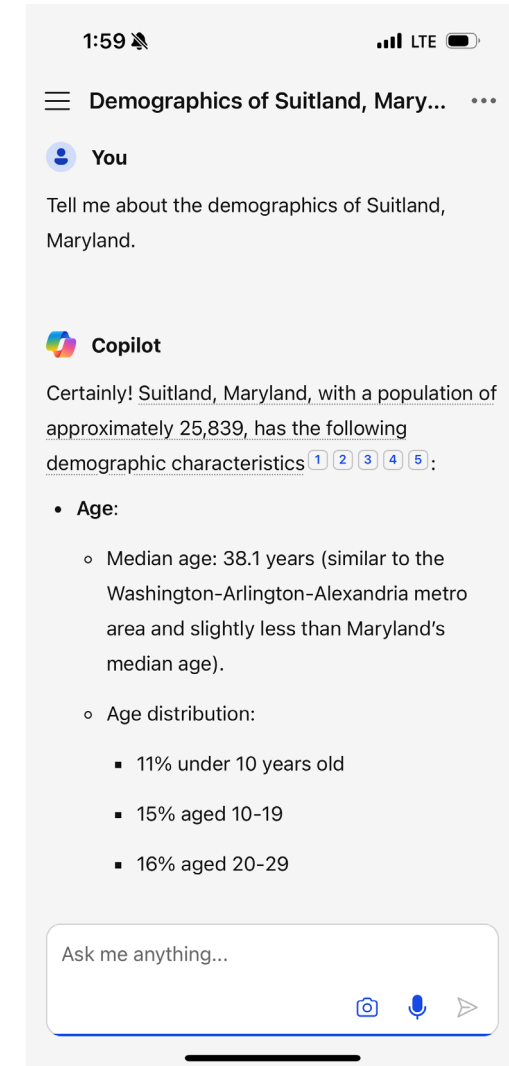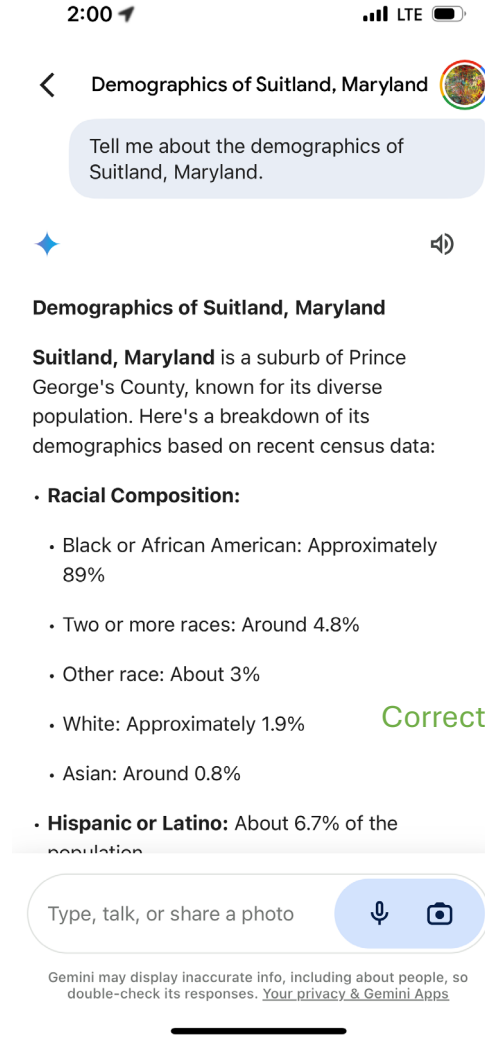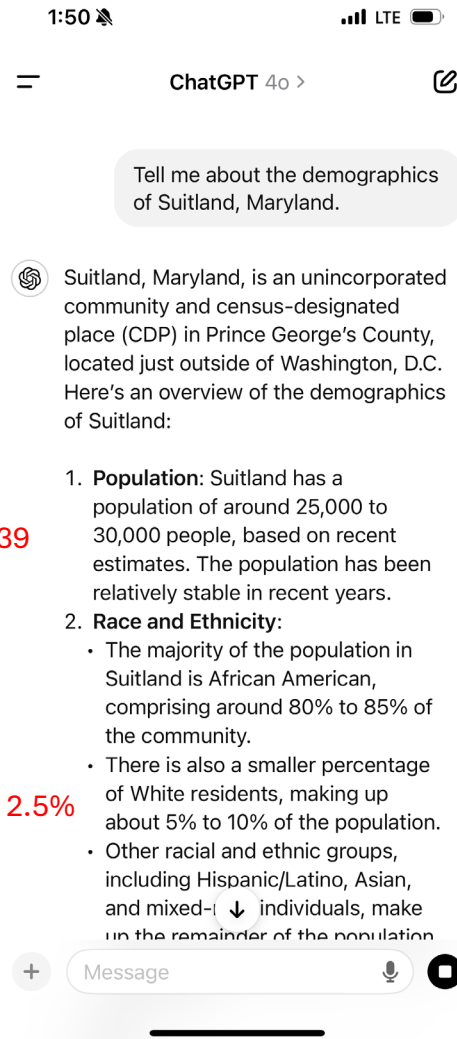
# ChatGPT returns incorrect results. Copilot cites secondary sources. Gemini shows improved accuracy. (June 2024)



**Left screenshot (ChatGPT):**

6:08

ChatGPT 4o

Tell me about the demographics of Suitland, Maryland.

Suitland, Maryland, located in Prince George's County, is a diverse suburb of Washington, D.C. The demographics of Suitland can be summarized as follows:

**Population**

- **Total Population:** Approximately 25,000 residents

*25,839*

**Racial and Ethnic Composition**

- **African American:** The majority, with about 85-90% of the population

*87.4%*

- **White:** Around 5-7%

*2.5%*

- **Hispanic or Latino:** Approximately 3-5%

*7.8%*

- **Asian:** About 1-2%
- **Other Races:** A small percentage includes mixed race and other racial groups

**Age Distribut...**

Message

**Middle screenshot (Gemini):**

6:09
◀ Safari

Suitland Demographics

Tell me about the demographics of Suitland, Maryland.

Suitland, Maryland has a predominantly Black or African American population, making up nearly 90% of residents. The median age in Suitland is around 38 years old, with a household income averaging $89,795. While the population has seen a slight increase in recent years, it remains a close-knit community with a rich history.

**Correct!**

Type, talk, or share a photo

Gemini may display inaccurate info, including about people, so double-check its responses. Your privacy & Gemini Apps

**Right screenshot (Copilot):**

6:12

Demographics of Suitland, ...

You

Tell me about the demographics of Suitland, Maryland.

Copilot

Explore

Certainly! Here's a summary of the demographics for Suitland, Maryland:

1: neighborhoodscount.com
2: censusreporter.org

- **Population**: As of 2020, Suitland has a population of approximately **25.84K** people [1].

- **Median Age**: The median age in Suitland is

Ask me anything...

United States® Census Bureau

# ChatGPT returns incorrect results. Copilot cites secondary sources. Gemini is accurate. (August 2024)



**ChatGPT 4o**

Tell me about the demographics of Suitland, Maryland.

Suitland, Maryland, is an unincorporated community and census-designated place (CDP) in Prince George's County, located just outside of Washington, D.C. Here's an overview of the demographics of Suitland:

1. **Population**: Suitland has a population of around 25,000 to 30,000 people, based on recent estimates. The population has been relatively stable in recent years.

2. **Race and Ethnicity**:
   - The majority of the population in Suitland is African American, comprising around 80% to 85% of the community.
   - There is also a smaller percentage of White residents, making up about 5% to 10% of the population.
   - Other racial and ethnic groups, including Hispanic/Latino, Asian, and mixed-race individuals, make up the remainder of the population

*25,839* (annotation)

*87.4%* (annotation)

*2.5%* (annotation)

---

**Demographics of Suitland, Maryland**

Tell me about the demographics of Suitland, Maryland.

## Demographics of Suitland, Maryland

**Suitland, Maryland** is a suburb of Prince George's County, known for its diverse population. Here's a breakdown of its demographics based on recent census data:

- **Racial Composition:**
  - Black or African American: Approximately 89%
  - Two or more races: Around 4.8%
  - Other race: About 3%
  - White: Approximately 1.9%
  - Asian: Around 0.8%
- **Hispanic or Latino:** About 6.7% of the population

*Correct!* (annotation)

Gemini may display inaccurate info, including about people, so double-check its responses. Your privacy & Gemini Apps

---

**Demographics of Suitland, Mary...**

**You**

Tell me about the demographics of Suitland, Maryland.

**Copilot**

Certainly! Suitland, Maryland, with a population of approximately 25,839, has the following demographic characteristics [1][2][3][4][5]:

- **Age:**
  - Median age: 38.1 years (similar to the Washington-Arlington-Alexandria metro area and slightly less than Maryland's median age).
  - Age distribution:
    - 11% under 10 years old
    - 15% aged 10-19
    - 16% aged 20-29

Ask me anything...

1: bing.com
2: censusreporter.org
3: areavibes.com
4: point2homes.com
5: datausa.io

*11.4%* (annotation)
*10.2%* (annotation)
*14.9%* (annotation)

United States® Census Bureau

Emerging Publishing Tactics

# Tactic #1: SEO, but More

# Data Discoverability:
# Sitemaps, Indexing, robots.txt, Server Side Rendering…

# What Is This Data About?
# Meta Tags, Open Graph Tags, schema.org, Croissant...

# Tell It To Me Straight:
# Documentation and Prose Descriptions
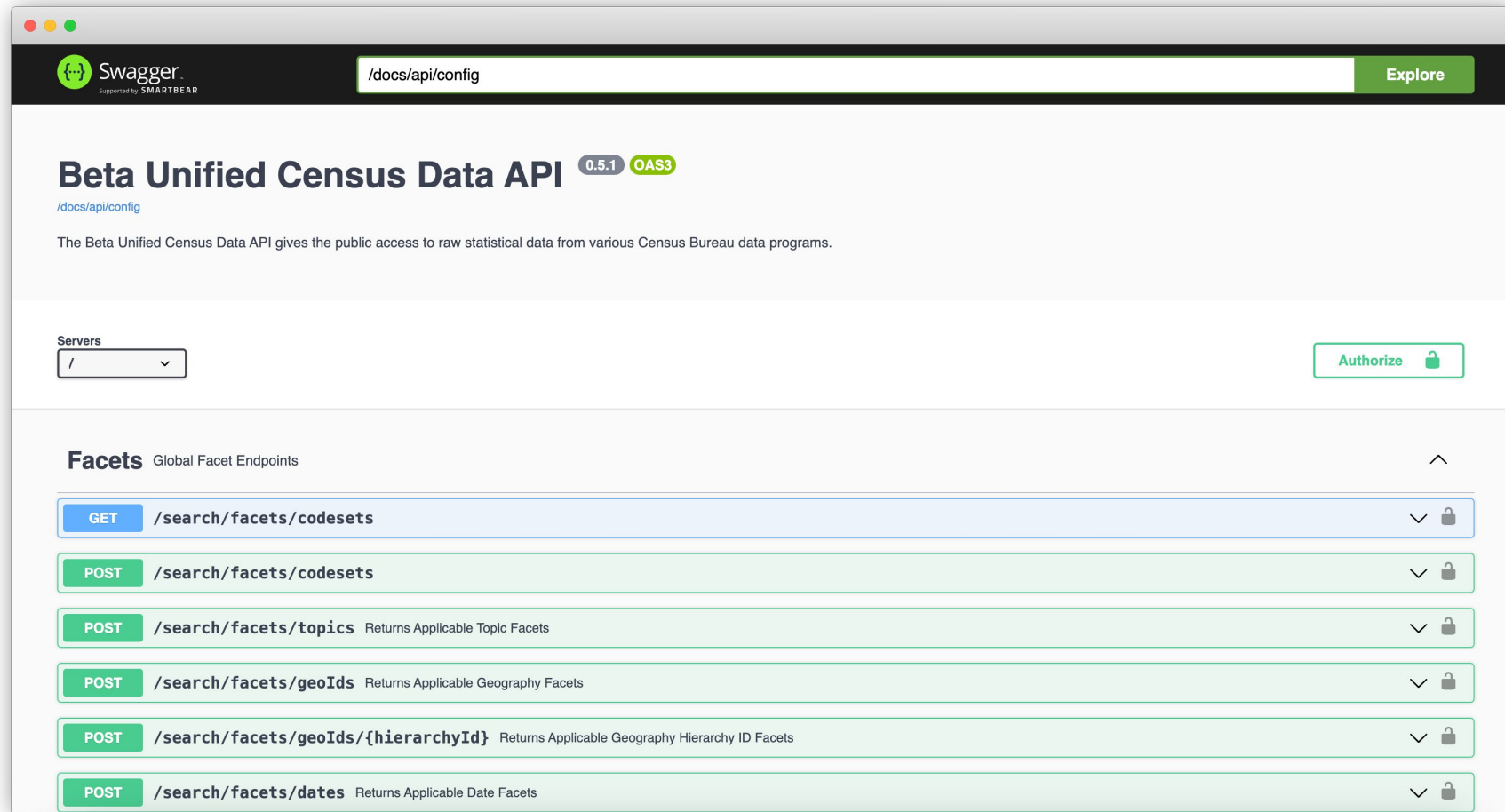
# Tactic #2: AI-Ready Data

United States® Census Bureau

# Aiding Data Ingestion:
# Proprietary and Unstructured < Open and Structured

XLSX  PDF     VS.     XML    CSV

                              JSON   Parquet

# Data Retrieval via API:
# Open Standards, Accessible Metadata

# Tabular Data:
# Hierarchical Tables < Flat Tables
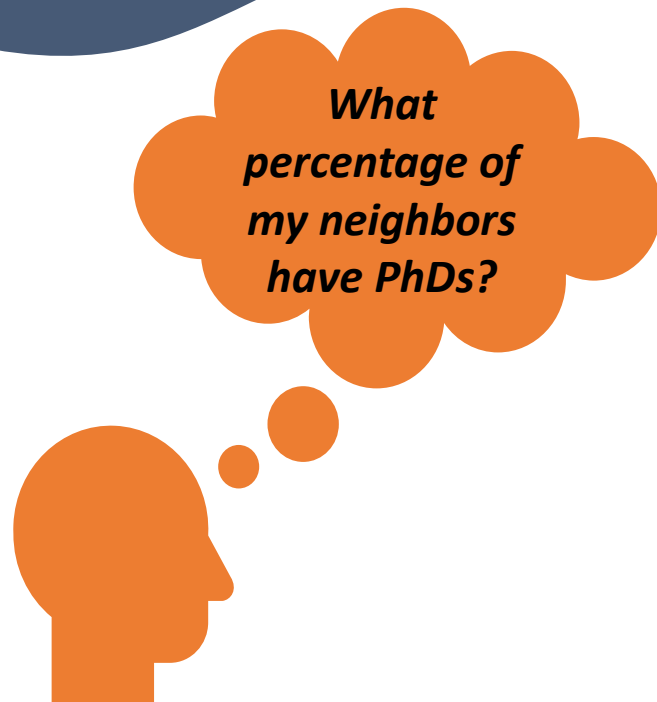
# Tactic #3: Meet AI Systems Where They Are

# Census API Assistant GPT
# (Research Prototype)

# Hybrid Solutions

Question

*What percentage of my neighbors have PhDs?*

LLM

Data Service

DB/API Query

Structured Response

Well Composed Answer
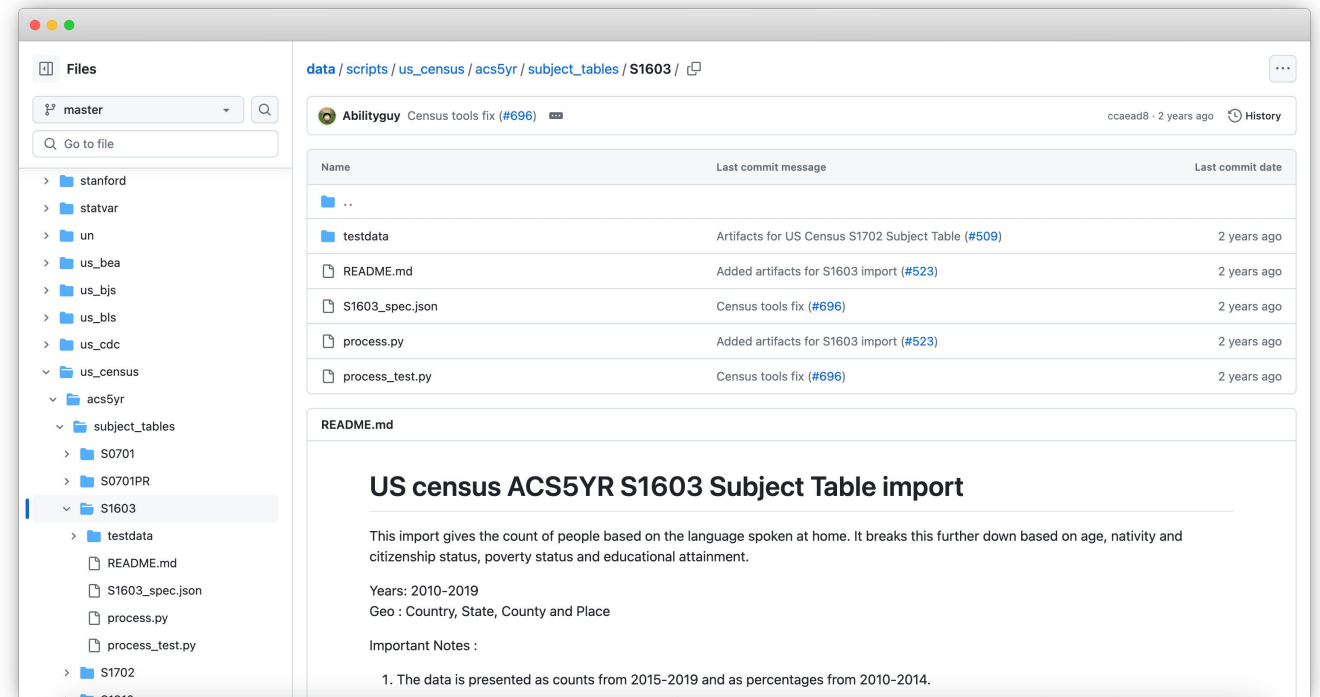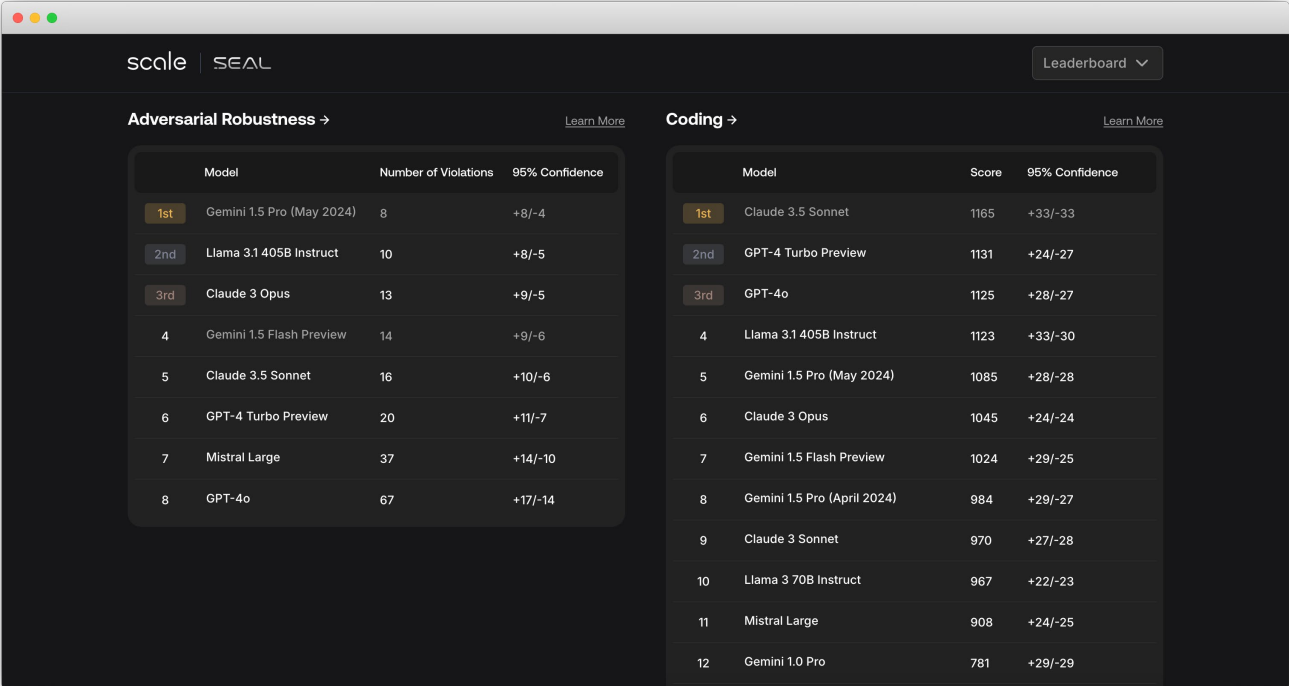
# Partnership and Engagement: Community

# Partnership and Engagement: Open Source Collaboration

# Partnership and Engagement: Benchmarks

scale | SEAL
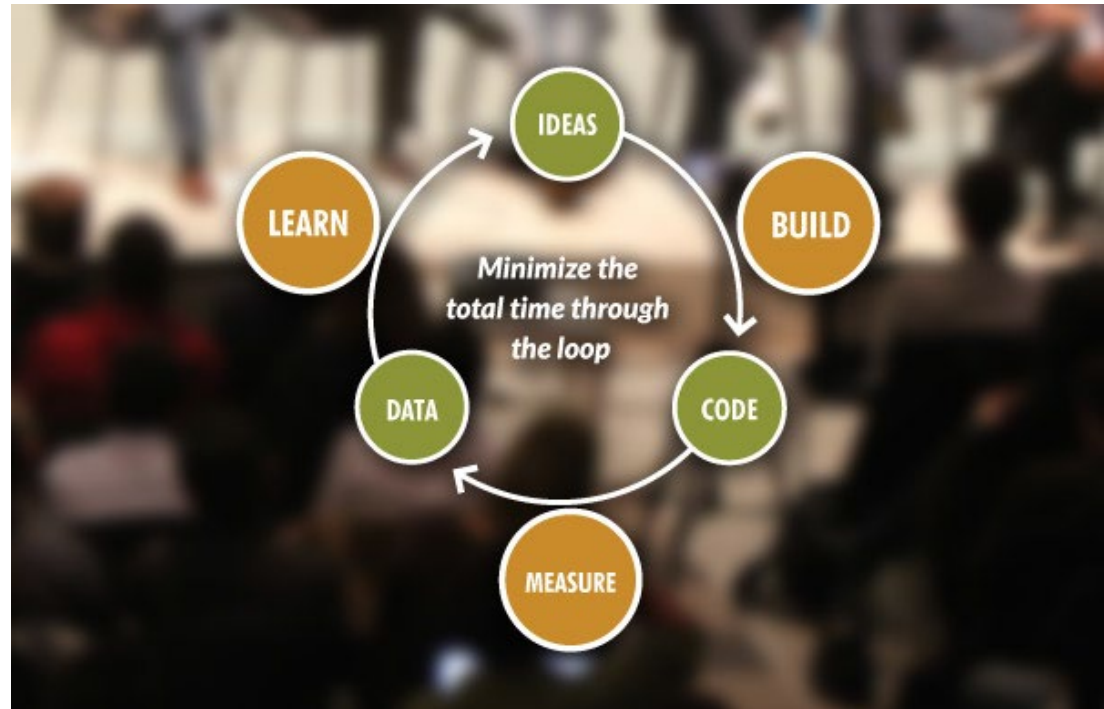
Leaderboard ⌄

**Adversarial Robustness →**  Learn More

| | Model | Number of Violations | 95% Confidence |
|---|---|---|---|
| 1st | Gemini 1.5 Pro (May 2024) | 8 | +8/-4 |
| 2nd | Llama 3.1 405B Instruct | 10 | +8/-5 |
| 3rd | Claude 3 Opus | 13 | +9/-5 |
| 4 | Gemini 1.5 Flash Preview | 14 | +9/-6 |
| 5 | Claude 3.5 Sonnet | 16 | +10/-6 |
| 6 | GPT-4 Turbo Preview | 20 | +11/-7 |
| 7 | Mistral Large | 37 | +14/-10 |
| 8 | GPT-4o | 67 | +17/-14 |

**Coding →**  Learn More

| | Model | Score | 95% Confidence |
|---|---|---|---|
| 1st | Claude 3.5 Sonnet | 1165 | +33/-33 |
| 2nd | GPT-4 Turbo Preview | 1131 | +24/-27 |
| 3rd | GPT-4o | 1125 | +28/-27 |
| 4 | Llama 3.1 405B Instruct | 1123 | +33/-30 |
| 5 | Gemini 1.5 Pro (May 2024) | 1085 | +28/-28 |
| 6 | Claude 3 Opus | 1045 | +24/-24 |
| 7 | Gemini 1.5 Flash Preview | 1024 | +29/-25 |
| 8 | Gemini 1.5 Pro (April 2024) | 984 | +29/-27 |
| 9 | Claude 3 Sonnet | 970 | +27/-28 |
| 10 | Llama 3 70B Instruct | 967 | +22/-23 |
| 11 | Mistral Large | 908 | +24/-25 |
| 12 | Gemini 1.0 Pro | 781 | +29/-29 |

United States® Census Bureau

What's Next?

United States® Census Bureau

# Monitor and Continue to Innovate (Cautiously and Ethically)

# THANK YOU

**Ryan Harper**, Product Manager
Center for Enterprise Dissemination
U.S. Census Bureau

ryan.harper@census.gov