# Prof. David Evans

## Data Privacy Expert

# Assignment

**1** Evaluate **privacy risks** with the proposed sharing of User-side Data, Ads Data, and Search Data

**2** Assess whether **privacy-enhancing technologies** can mitigate those privacy risks while still sharing useful information

**3** Respond to the reports of Google's privacy expert

REDACTED FOR PUBLIC FILING

# Key Opinions

**1** There are well-established **privacy-enhancing techniques** that can be used to protect sensitive information.

**2** Many organizations, including Google, **safely release sensitive data** by using privacy-enhancing techniques.

**3** Google can share the data at issue in a way that **assures privacy while providing utility**.

# Google's Expert Agrees Data Can Be Shared



**Chris Culnane, PhD**

**Google's Expert**
**Principal & Consultant**
**Castellate Consulting Ltd.**

Q. Dr. Culnane, you believe that it is possible for Google to share what you call the DOJ search data by applying privacy-enhancing techniques to achieve suitable privacy safeguards, don't you?

**A. Yes.**

# Experts' Disagreement

## What Dr. Culnane Claims

"In the Search Context, Only Frequency Thresholds Provide Indistinguishability."
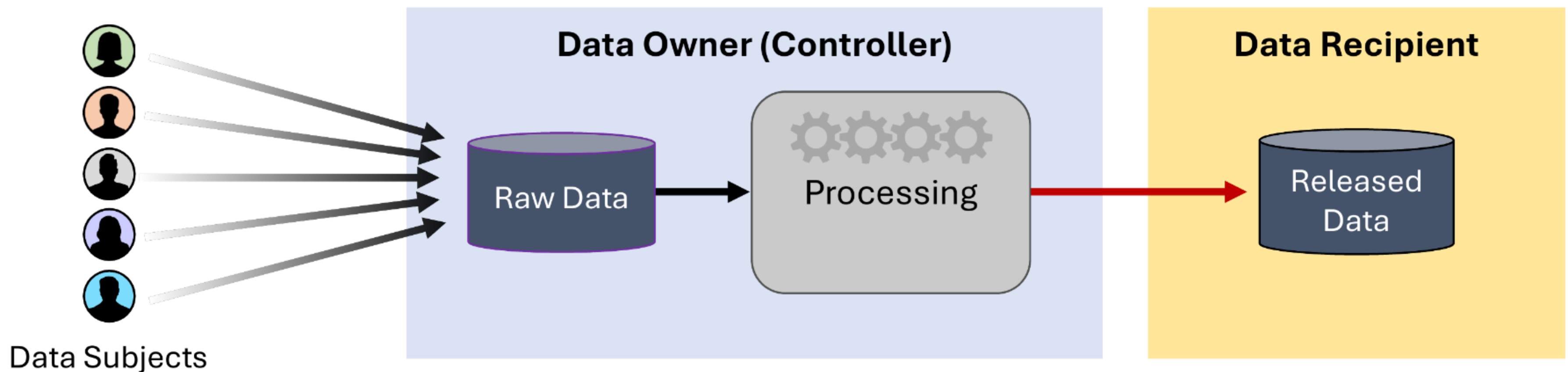
## My Opinion

There are many well-established **privacy-enhancing techniques**, and the remedy should **use techniques appropriately to assure privacy while providing high utility.**

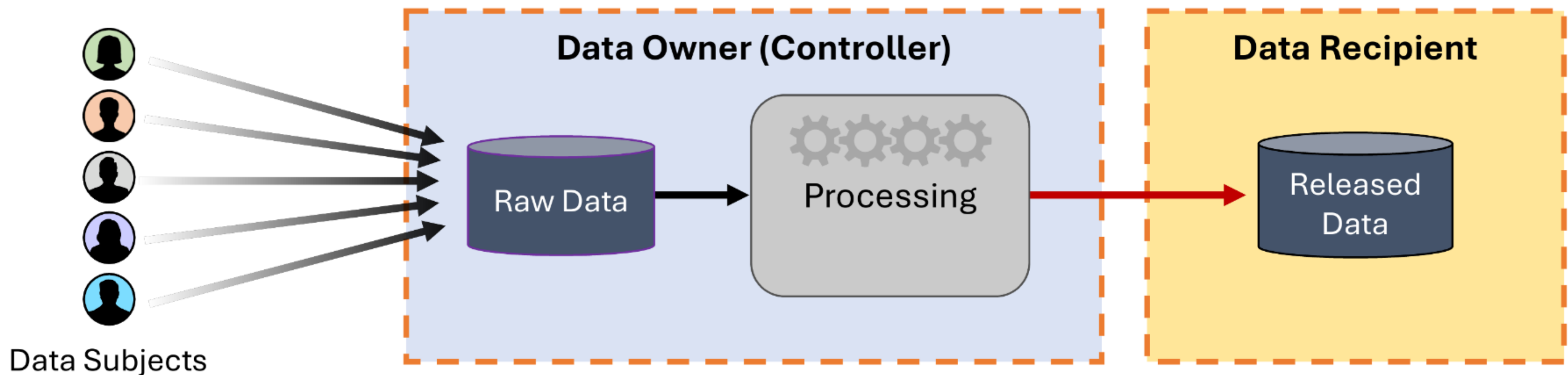# What is Data Privacy?



**Data Collection**

**Data Processing**

**Data Release**

Data Owner (Controller)

Data Recipient

Data Subjects

Raw Data

Processing

Released Data

# What is Data Privacy?

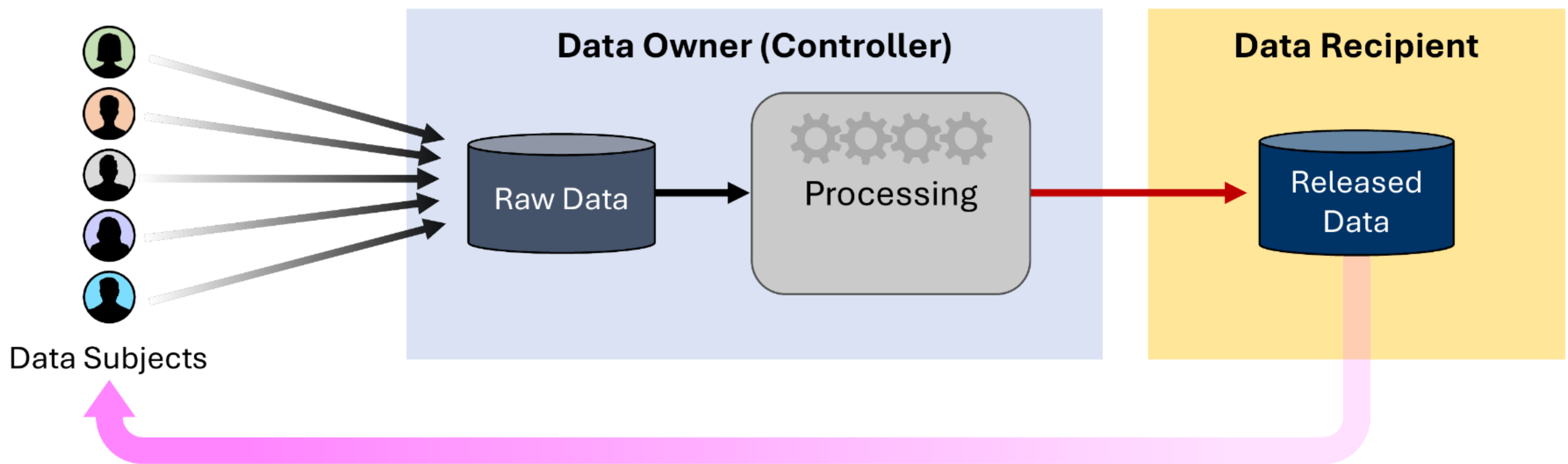**Data Collection**　　　　　**Data Processing**　　　　　**Data Release**



**Data security**: preventing **unintended releases** of data

# What is Data Privacy?

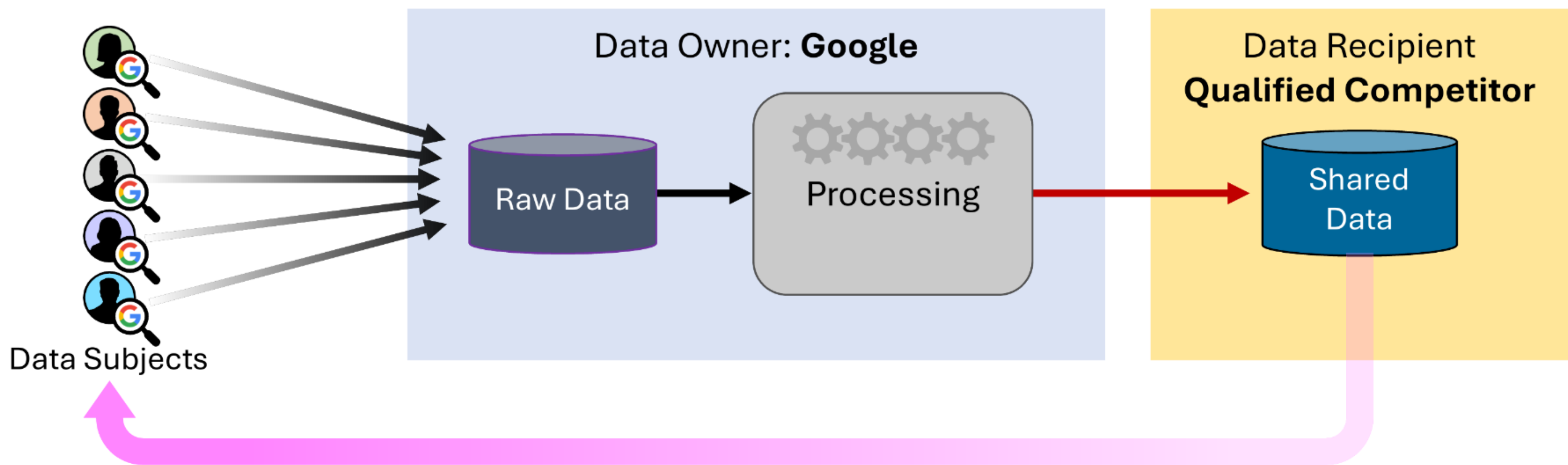**Data Collection**  **Data Processing**  **Data Release**



**Data privacy:** preventing **unintended disclosure** of sensitive information from **intentionally released data**

# Data Privacy for Proposed Data Sharing

**Data Collection**

**Data Processing**

**Data Release**

Data Subjects

Data Owner: **Google**

Raw Data → Processing

Data Recipient **Qualified Competitor**

Shared Data

**Data privacy issue:** potential for **disclosure** of sensitive information from **shared data** and mitigations to share safely

# The Data at Issue

**User-side Data**
RPFJ Sections VI.A, C, & D

**Search Index Data**
RPFJ Section VI.A

**Ads Data**
RPFJ Sections VI.E & F



Submitted queries
Clicked-on links
Time looking at results
Hovering over a link
User location
User device
Ranking signals
...

**Data Google collects from users and uses to train models**
**(RankEmbed, NavBoost, Glue, and** ⬚ **)**

# Innocuous Data Can Reveal Sensitive Information


**Linking**


Aggregate Statistics

| Block | Total | Race 1 | ... | Race 63 | ... |
|-------|-------|--------|-----|---------|-----|
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |
| 20394 | 712 | 0 | ... | 82 | ... |
| 20395 | 2316 | 3 | ... | 27 | ... |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |

| ID | Block | Race | Age | ... |
|-----|-------|------|-----|-----|
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |
| 032 | 20394 | 7 | 23 | ... |
| 033 | 20394 | 5 | 82 | ... |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |

Individuals

**Reconstruction**


Attribute Inference Attacks

Membership Inference Attacks
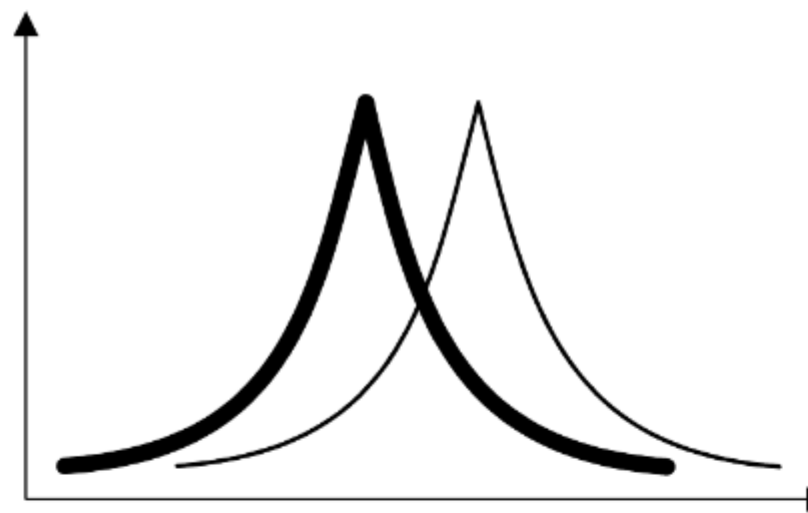
**Inference Attacks**

# Assessing Privacy Risk

## Until ~2000:
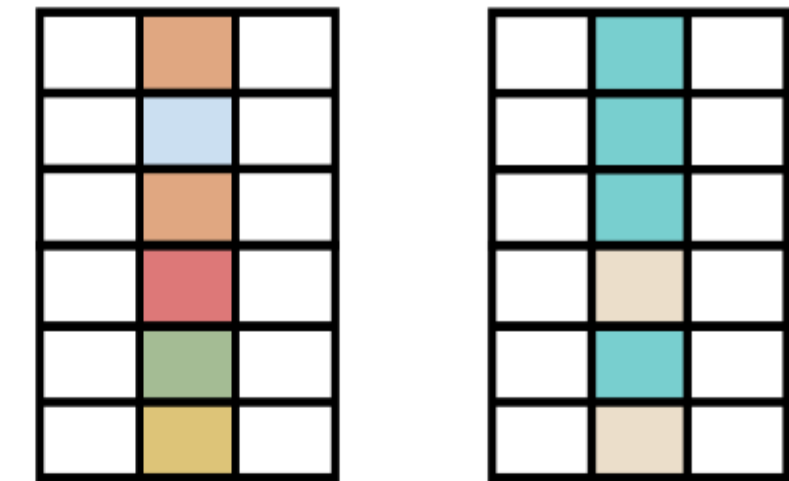## ad hoc privacy

Trying things and
hope they work

## Today: formal privacy

Mathematical definitions of privacy and
principled mechanisms for satisfying them
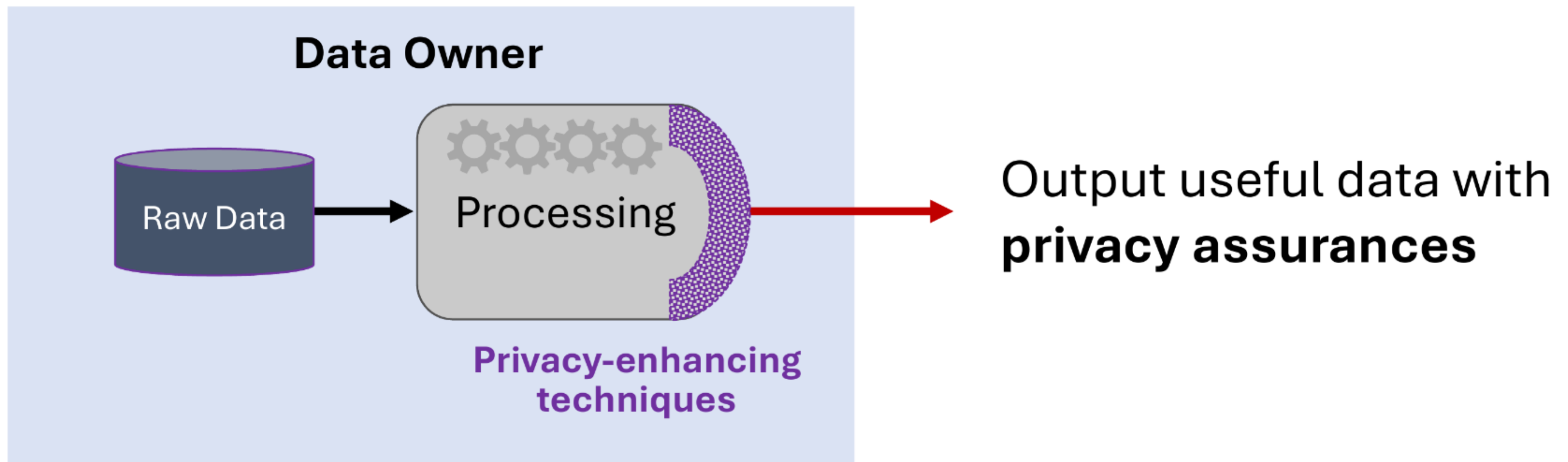


**differential privacy**

**_k_-anonymity**

and hundreds of others...

# Privacy-Enhancing Techniques (PETs)

**Data Owner**

Raw Data → Processing

**Privacy-enhancing techniques**

Output useful data with **privacy assurances**

# Opinion 1: Privacy-Enhancing Techniques Work

**1** There are well-established **privacy-enhancing techniques** that can be used to protect sensitive information.

**2** Many organizations, including Google, **safely release sensitive data** by using privacy-enhancing techniques.
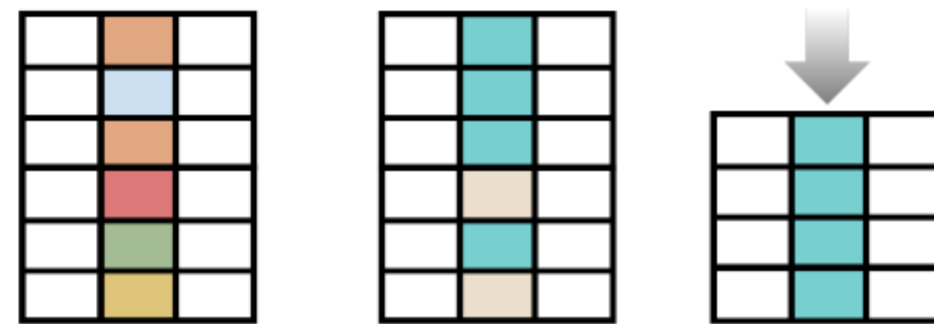
**3** Google can share the data at issue in a way that **assures privacy** while **providing utility**.

# Broad Types of Privacy-Enhancing Techniques



**Noise**

**Frequency Bounds**

**Cryptographic Methods**

Google

Visa

Private Set Intersection

Encrypted Inputs

Ad views connected with offline purchases

Source: Ion, Mihaela, Kreuter, Ben et al. "On Deploying Secure Computing: Private Intersection-Sum-with-Cardinality". In: 2020 IEEE European Symposium on Security and Privacy (EuroS&P), 2020, pp. 370–389, url: https://ieeexplore.ieee.org/document/9230369.

# Noise for Privacy

**Source Data**     **Noise**     **Released Data**

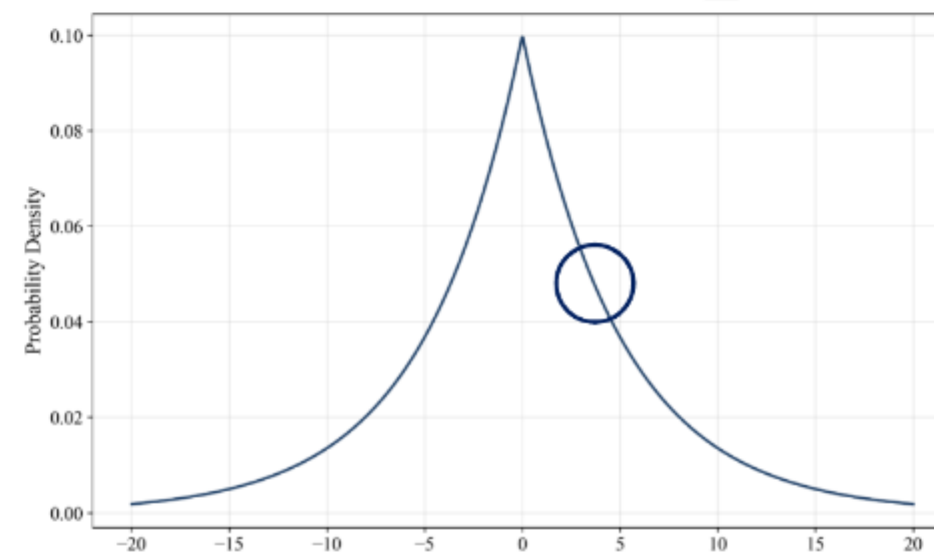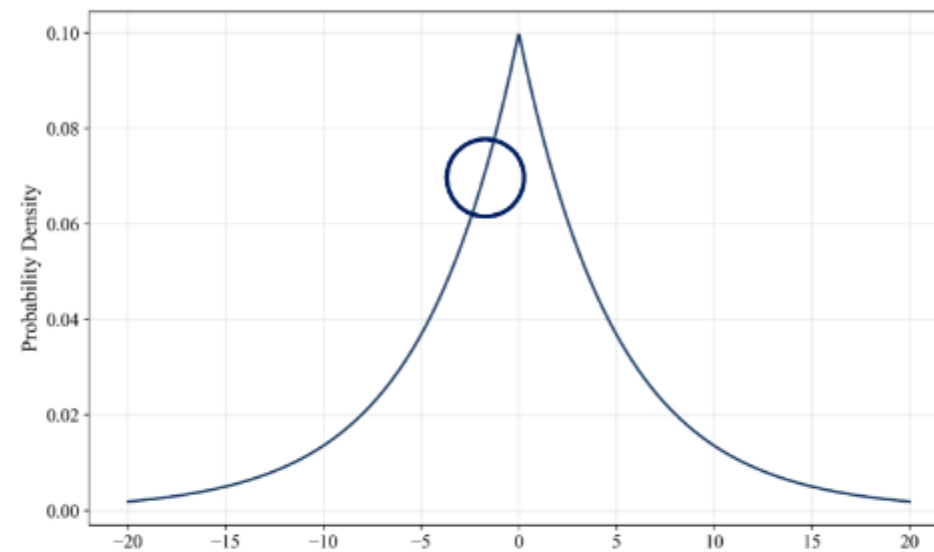629     +          =     631.52

629     +          =     628.73

# Differential Privacy

Gives a **mathematical bound** on exposure of individual's data

**No assumptions needed** about what is sensitive information, actual data, what adversary can do, what adversary already knows

$$\frac{\text{Probability of this output from dataset containing user}}{\text{Probability of this output from dataset \textbf{without} user}} \leq \exp(\epsilon)$$
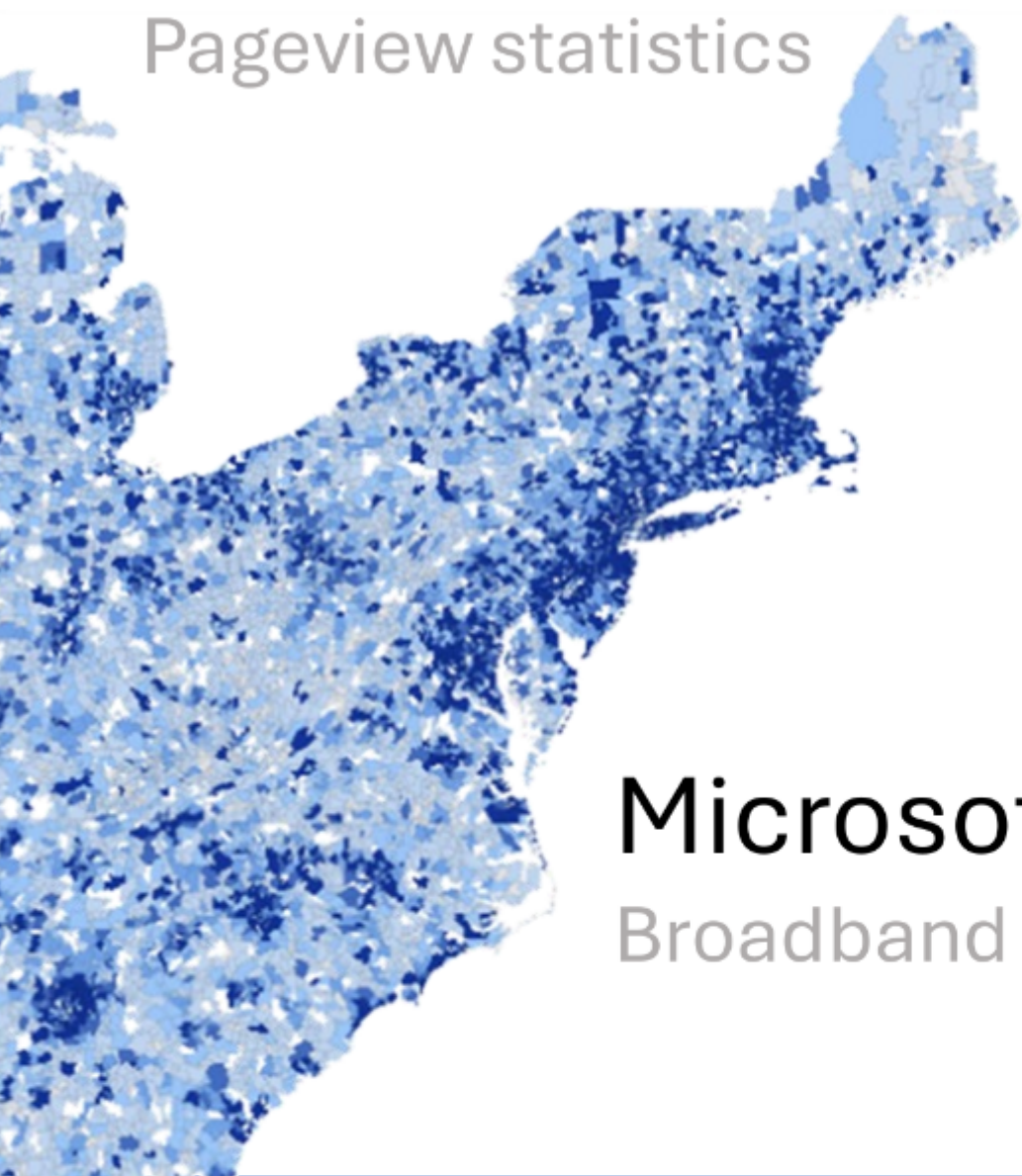
Privacy loss parameter (epsilon) provides precise control of **privacy–utility tradeoff**

# Widespread Acceptance and Use

**Wikimedia**
Pageview statistics

**Apple**
Learning iconic scenes

Solano   Sonoma

**Facebook**
Movement dataset

**Microsoft**
Broadband usage

**US Census Bureau**
Redistricting data

US Census

no formal privacy guarantee

ε-differential privacy

1980  1990  2000  2010  2020

# Google Uses Differential Privacy (DP)



**Internal Google Document**



**Google AI Comic**



**Variation in Mobility**

REDACTED FOR PUBLIC FILING

# *K*-anonymity Formal Privacy Definition

Privacy definition that requires that any released data record is **indistinguishable** from at least $k - 1$ other records.

| Query | Location | Device |
|---|---|---|
| best Mexican food | (38.8977°, 77.0365°) | Pixel9a-Android15-v23.523 |

**Indistinguishable**

| Query | Location | Device |
|---|---|---|
| best Mexican food | (38.8977°, 77.0365°) | Pixel9a-Android15-v23.523 |

**Indistinguishable?**

| Query | Location | Device |
|---|---|---|
| best **m**exican food | (38.8977°, 77.036**4**°) | Pixel9a-Android15-v**22.173** |

**Indistinguishable?**

| Query | Location | Device |
|---|---|---|
| mexican **restaurant** | (38.8977°, 77.0365°) | Pixel9a-Android15-v**21.083** |

# How to Satisfy *K*-anonymity

Source data
(*k*=**1**)

| Query | Location | Device |
|---|---|---|
| best Mexican food | (38.8977°, 77.0365°) | Pixel9a-Android15-v23.523 |
| best Mexican food | (38.8977°, 77.0365°) | Pixel9a-Android15-v23.523 |
| best **m**exican food | (38.8977°, 77.036**4**°) | Pixel9a-Android15-v**22.173** |
| mexican **restaurant** | (38.8977°, 77.0365°) | Pixel9a-Android15-v**21.083** |

Record removal ⬇

Released data
(*k*=**2**)

| Query | Location | Device |
|---|---|---|
| best Mexican food | (38.8977°, 77.0365°) | Pixel9a-Android15-v23.523 |
| best Mexican food | (38.8977°, 77.0365°) | Pixel9a-Android15-v23.523 |
| best **m**exican food | (38.8977°, 77.036**4**°) | Pixel9a-Android15-v**22.173** |
| mexican **restaurant** | (38.8977°, 77.0365°) | Pixel9a-Android15-v**21.083** |

# How to Satisfy *K*-anonymity with Utility

Source data
(*k*=**1**)

| Query | Location | Device |
|---|---|---|
| best Mexican food | (38.8977°, 77.0365°) | Pixel9a-Android15-v23.523 |
| best Mexican food | (38.8977°, 77.0365°) | Pixel9a-Android15-v23.523 |
| best **m**exican food | (38.8977°, 77.036**4**°) | Pixel9a-Android15-v**22.173** |
| mexican **restaurant** | (38.8977°, 77.0365°) | Pixel9a-Android15-v**21.083** |

Generalization ⬇        Suppression ⬇

Released data
(*k*=**3**)

| Query | Location | Device |
|---|---|---|
| best Mexican food | **DC 20500** | Pixel9a-Android15▓▓▓▓▓ |
| best Mexican food | **DC 20500** | Pixel9a-Android15▓▓▓▓▓ |
| best **M**exican food | **DC 20500** | Pixel9a-Android15▓▓▓▓▓ |
| mexican **restaurant** | DC 20500 | Pixel9a-Android15-v**21.083** |

# Better Generalization Improves Utility

**Source data (*k*=1)**

| Query | Location | Device |
|---|---|---|
| best Mexican food | (38.8977°, 77.0365°) | Pixel9a-Android15-v23.523 |
| best Mexican food | (38.8977°, 77.0365°) | Pixel9a-Android15-v23.523 |
| best **m**exican food | (38.8977°, 77.036**4**°) | Pixel9a-Android15-v**22.173** |
| mexican **restaurant** | (38.8977°, 77.0365°) | Pixel9a-Android15-v**21.083** |

Generalization ⬇    Suppression ⬇

**Released data (*k*=4)**

| Query Intent | Location | Device |
|---|---|---|
| Mexican restaurant | **DC 20500** | Pixel9a-Android15▨▨▨ |
| Mexican restaurant | **DC 20500** | Pixel9a-Android15▨▨▨ |
| Mexican restaurant | **DC 20500** | Pixel9a-Android15▨▨▨ |
| Mexican restaurant | **DC 20500** | Pixel9a-Android15▨▨▨ |

# Example uses Generalization for *K*-anonymity

**CDC**

Public Use Data

**Generalization**
**Partial Suppression**
**L-diversity**

**Cloudflare**

Validating Leaked Passwords

**Generalization**
**Partial Suppression**

**Facebook**

URLs Dataset

**Generalization**
**Partial Suppression**
**Differential Privacy**

HARVARD UNIVERSITY

**SOCIAL SCIENCE ONE**
Hosted by Harvard's Institute for Quantitative Social Science

**RFP for URL Shares**

This is a codebook for data on the demographics of people who viewed, shared, and otherwise interacted with web pages (URLs) shared on Facebook. The data has about 68

# Uses of Generalization for Privacy at Google



**COVID-19 Vaccination Search Insights**
**Generalization** (Geographic, Time, Grouping search queries)



"[W]e use generalization to remove a portion of the data or replace some part of it with a common value. . . . Generalization allows us to achieve k-anonymity . . . ."

**Google's Privacy Policy**
**Generalization for k-anonymity**

# Google's Data Sharing Implementation For DMA

Anonymity Set A (Size 5)

| Query | Country | Device |
|---|---|---|
| News today | DE | Mobile |
| News today | DE | Desktop |
| News today | DE | Desktop |
| News today | DE | Mobile |
| News today | DE | Mobile |

Anonymity Set A.1 (Size 2)

| Query | Country | Device |
|---|---|---|
| News today | DE | Desktop |
| News today | DE | Desktop |

Anonymity Set A.2 (Size 3)

| Query | Country | Device |
|---|---|---|
| News today | DE | Mobile |
| News today | DE | Mobile |
| News today | DE | Mobile |

**No** Field Suppression
**No** Generalization
**No** Spell-Correcting Queries
**No** Grouping by Query Intent

**Google's Experts' Report on DMA**
(Dr. Culnane and Prof. Rubenstein)

21. Google identified three additional recovery mechanisms and is working on implementing them. These mechanisms require significant engineering work to develop and will therefore not be ready for the initial dataset, but Google expects to introduce them for the second quarterly release of its Art. 6(11) dataset.

22. First, Google has developed a privacy-safe way to release additional data about low-volume queries. For queries that typically fail to meet th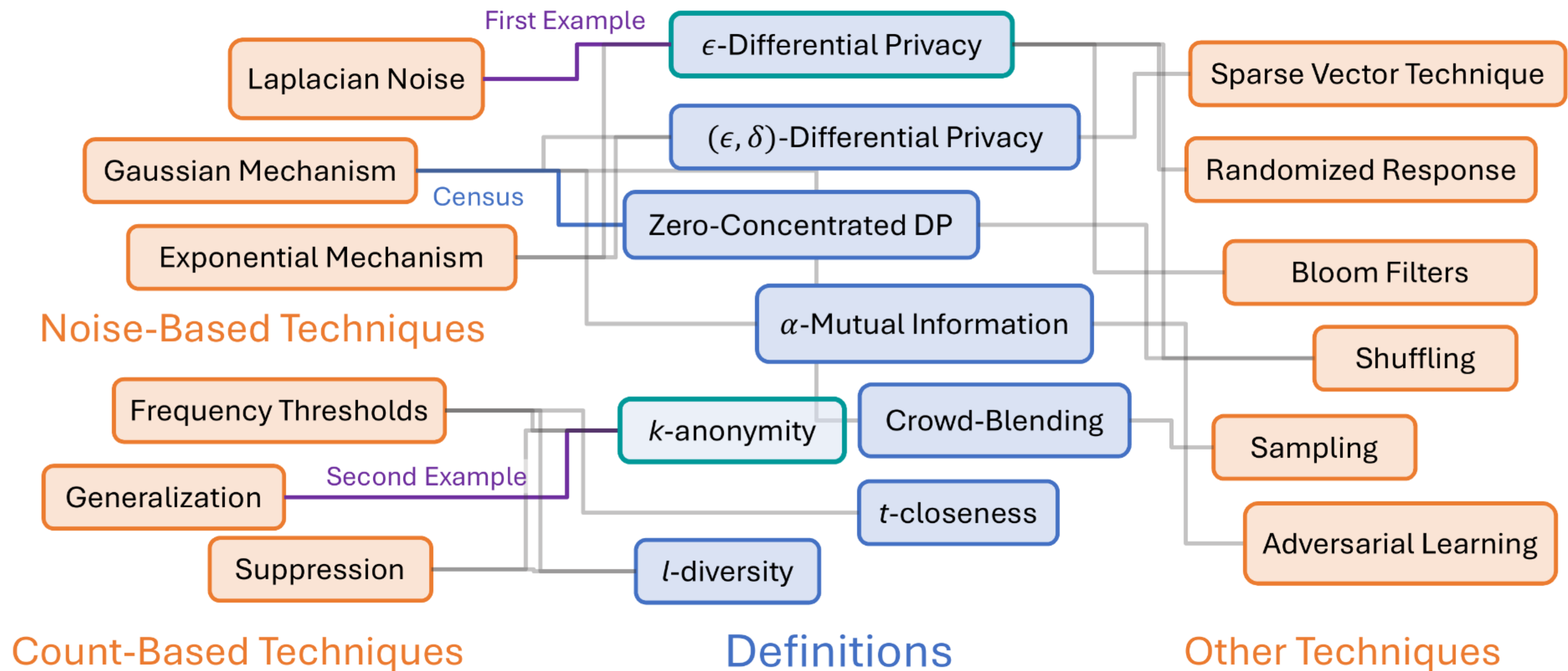e m-threshold for a given country, Google will apply the threshold at an EEA-wide level, and report combined statistics across the EEA instead of suppressing the data for many queries that do not support finer country-level data.

*Generalization by combining all countries*

23. Second, Google Search automatically corrects some typos and misspellings in user queries, showing the user results for the corrected query. Before anonymization, Google will replace "typo" queries that were automatically corrected for the results shown to the user with their corrected versions.

*Generalization by fixing "typo" queries*

24. Third, Google has developed an additional mechanism to "map" certain low-frequency queries that Search does not automatically correct (e.g., [mssql

**Google's Second Response to European Commission**
(January 2024, 1¼ years after DMA)

# Many Formal Privacy Definitions And Principled Techniques



Noise-Based Techniques

Laplacian Noise — First Example → $\epsilon$-Differential Privacy

Gaussian Mechanism — Census

Exponential Mechanism

$(\epsilon, \delta)$-Differential Privacy

Zero-Concentrated DP

$\alpha$-Mutual Information

Frequency Thresholds

Generalization — Second Example → $k$-anonymity

Suppression

Crowd-Blending

$t$-closeness

$l$-diversity

Count-Based Techniques

Definitions

Sparse Vector Technique

Randomized Response

Bloom Filters

Shuffling

Sampling

Adversarial Learning

Other Techniques

# Opinion 2: PETs Can Be Used To Safely Release Useful Data

**1**   There are well-established **privacy-enhancing techniques** that can be used to protect sensitive information.

**2**   Many organizations, including Google, **safely release sensitive data** by using privacy-enhancing techniques.

**3**   Google can share the data at issue in a way that **assures privacy** while **providing utility.**

**REDACTED FOR PUBLIC FILING**

# Selecting Appropriate Privacy-Enhancing Techniques

## Properties of the **source data**

- Type and amount
- Granularity
- Dimensionality
- Sensitivity
- Update frequency

...

### Disclosure Risk

## Uses of the **released data**

- Amount required
- Granularity needed
- Correlations used
- Accuracy thresholds
- Sharing frequency

...

### Data Utility

# Selecting Privacy-Enhancing Techniques for Data at Issue

**Slide 37 from Google's Opening Statement**
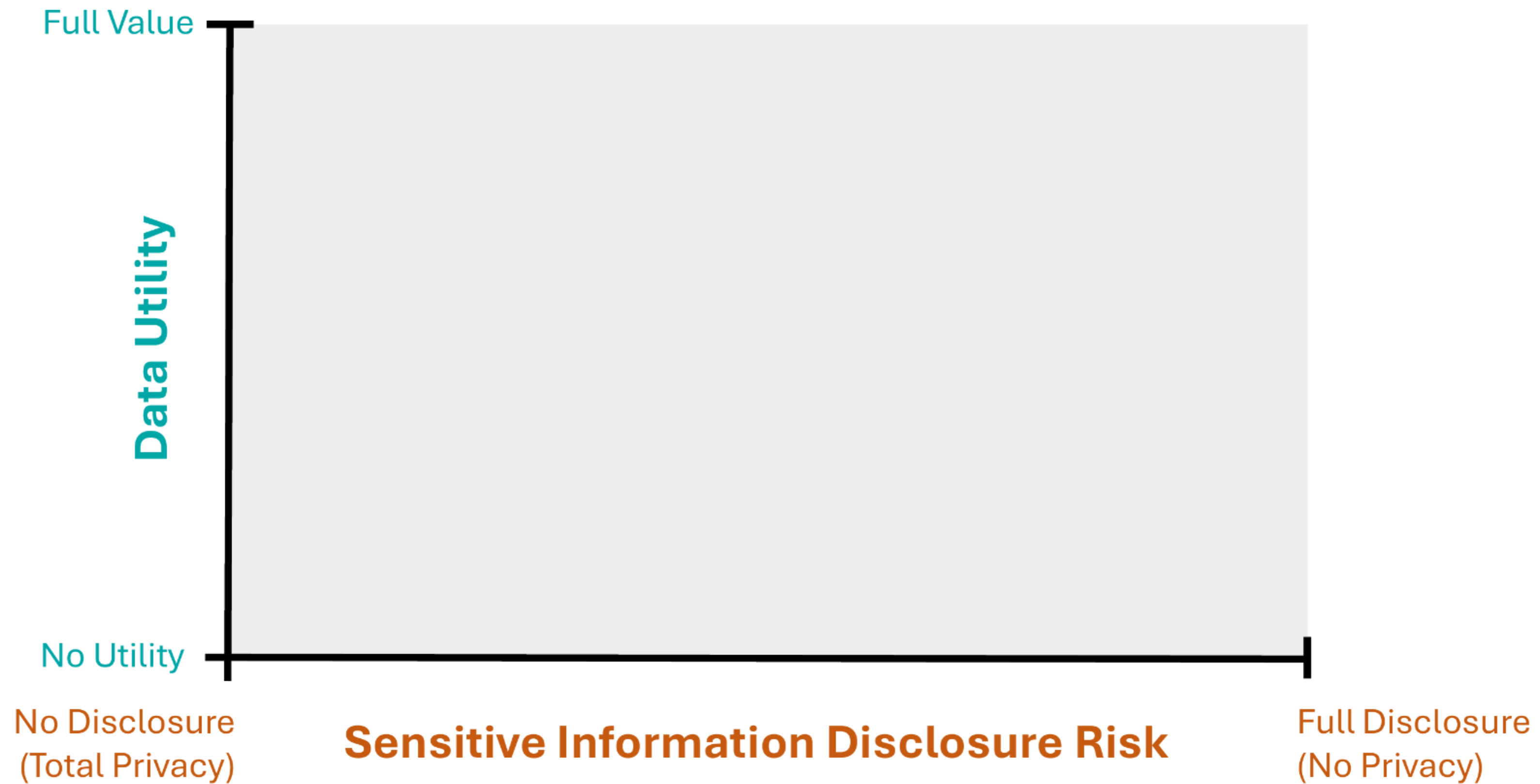
## Plaintiffs' Privacy Expert Offers No Opinion

**David Evans, PhD.**
DOJ Expert

A. There are many ways to protect text data, and one way is to use the frequency-based method to achieve a definition similar to K-anonymity.

Q. That is what you propose should be done here?

A. I don't make any proposal as to what should be done here. I just speak to the availability of many different privacy-enhancing techniques that could be used to satisfy the requirements of the RPFJ.
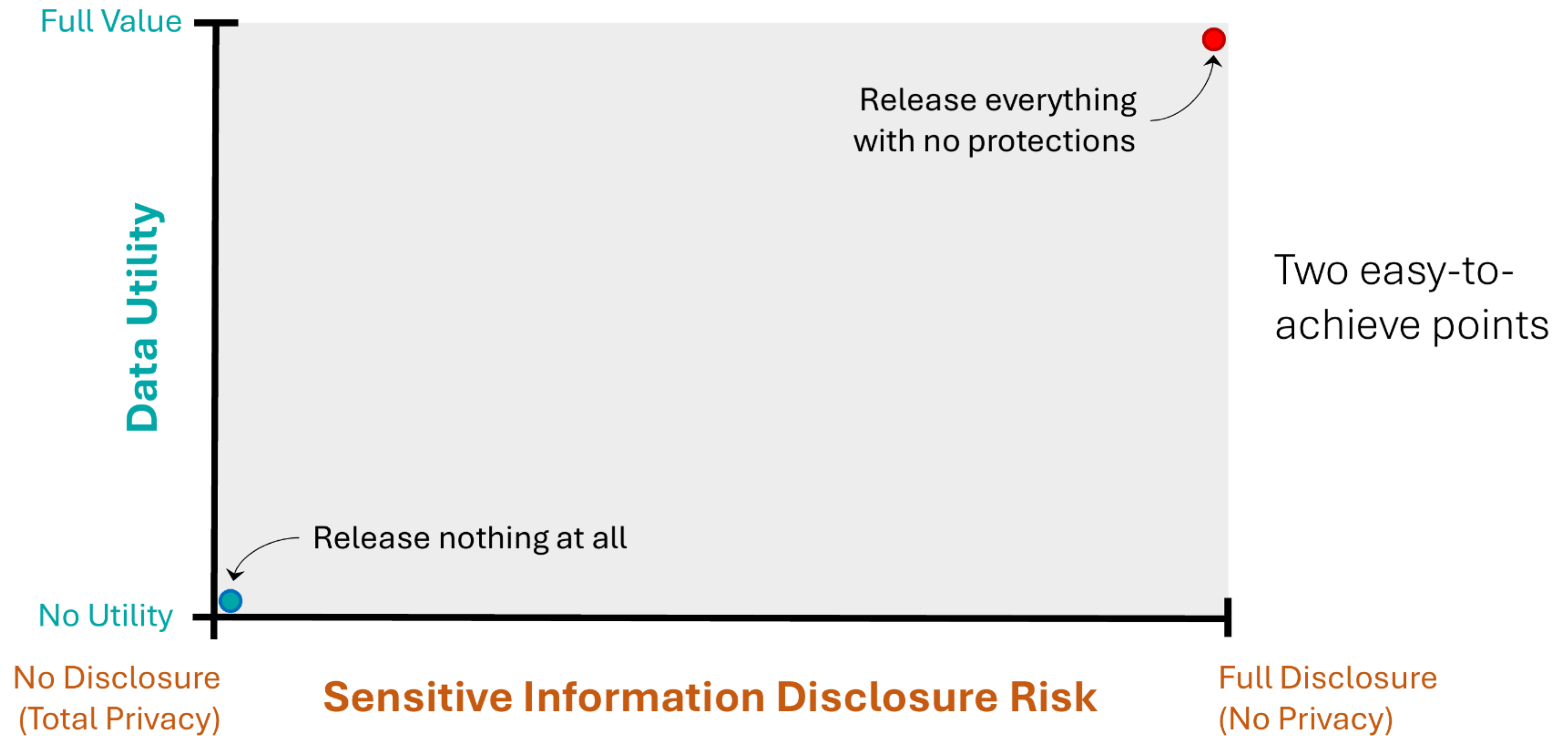
Evans (DOJ) Trial Tr. 130:10-22
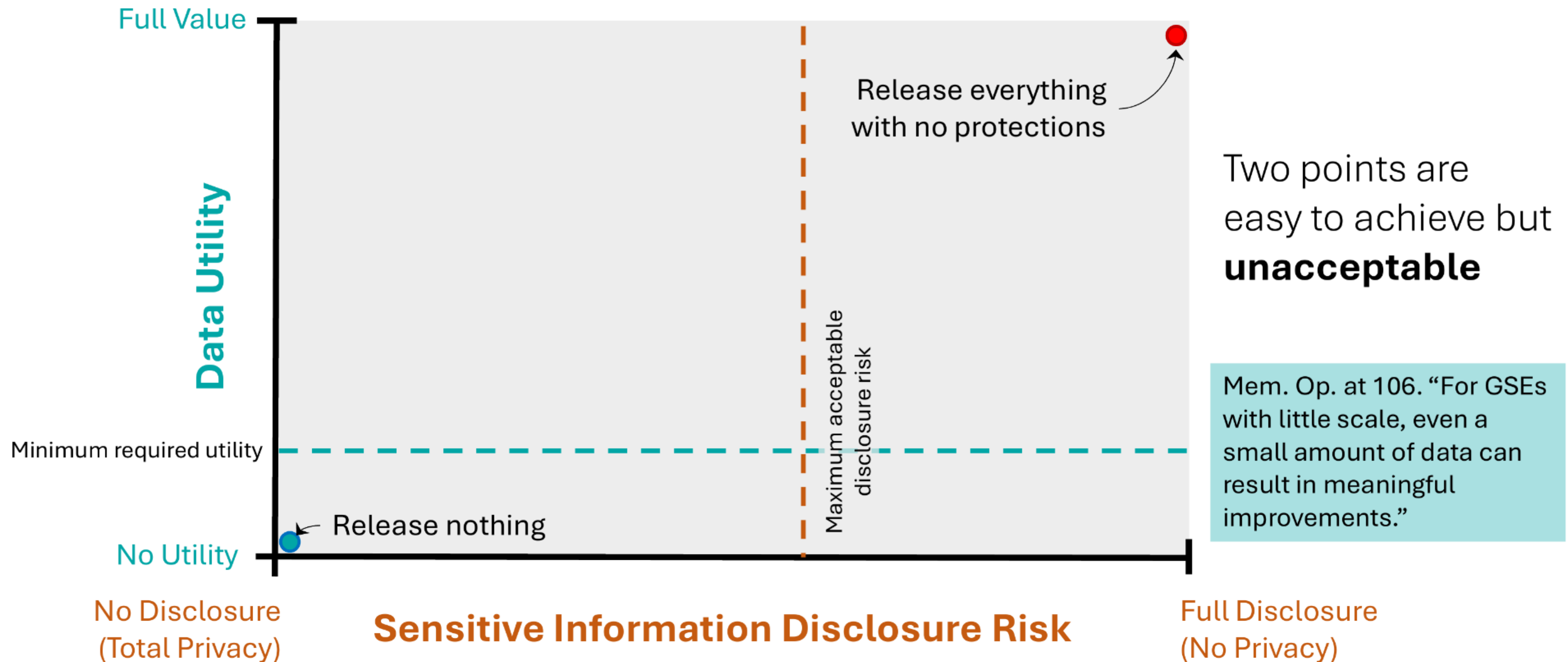
Google
RDXD-01.037

# Privacy–Utility Tradeoff

REDACTED FOR PUBLIC FILING

# Privacy–Utility Tradeoff



Source: Adapted from Opening Report Figure 1: Privacy – utility tradeoff curve.

# Privacy–Utility Tradeoff



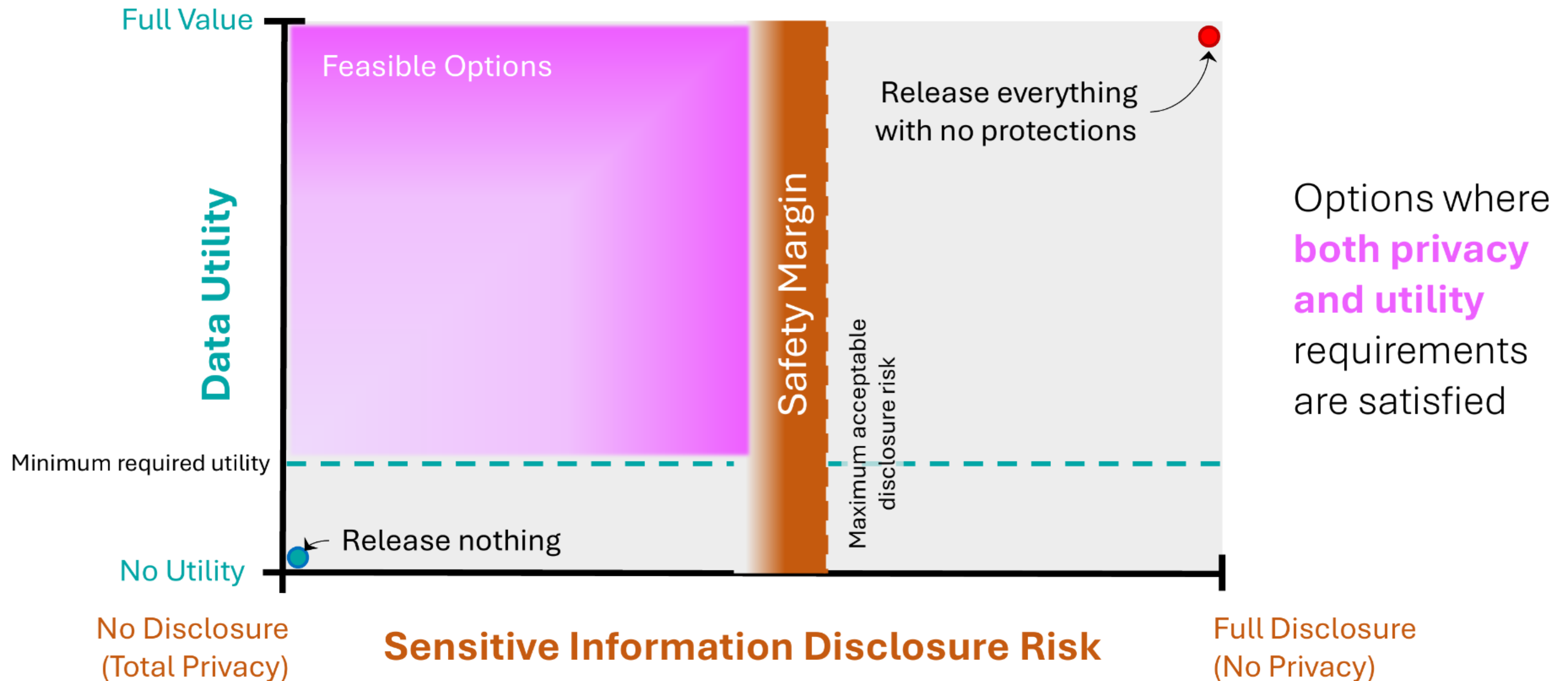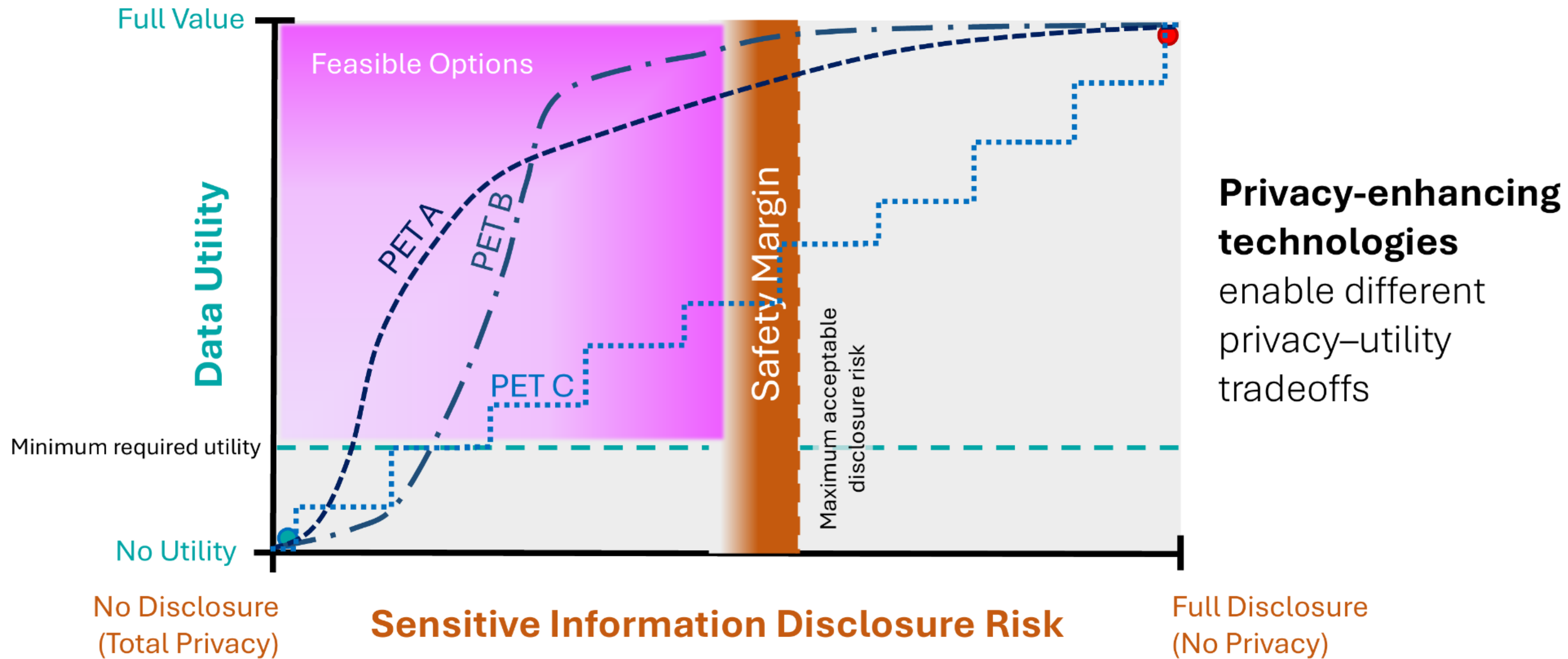Two points are easy to achieve but **unacceptable**

Mem. Op. at 106. "For GSEs with little scale, even a small amount of data can result in meaningful improvements."

# Privacy–Utility Tradeoff

# Privacy–Utility Tradeoff



**Privacy-enhancing technologies** enable different privacy–utility tradeoffs

Figure labels: Full Value, Feasible Options, Data Utility, Minimum required utility, No Utility, PET A, PET B, PET C, Safety Margin, Maximum acceptable disclosure risk, Sensitive Information Disclosure Risk, No Disclosure (Total Privacy), Full Disclosure (No Privacy)

# Privacy–Utility Tradeoff

Full Value

Feasible Options

Data Utility

PET A

PET B

PET C

Safety Margin

Minimum required utility

No Utility

PET D

No Disclosure
(Total Privacy)

**Sensitive Information Disclosure Risk**

Full Disclosure
(No Privacy)

Privacy-enhancing technologies could (**but never should**) be used in ways that **reduce both privacy and utility**

# Privacy–Utility Tradeoff



Combinations of techniques often provide the best privacy–utility tradeoff

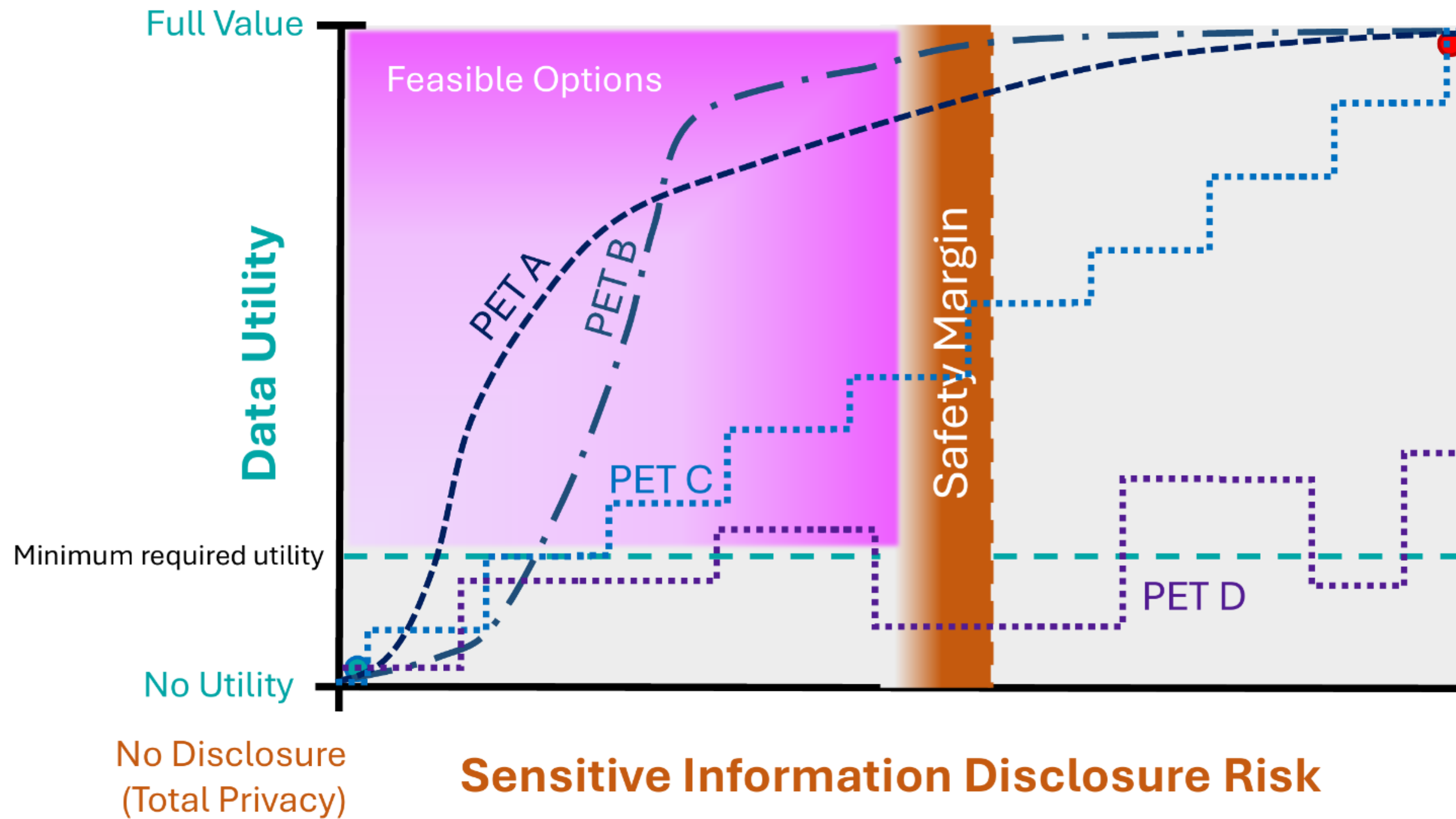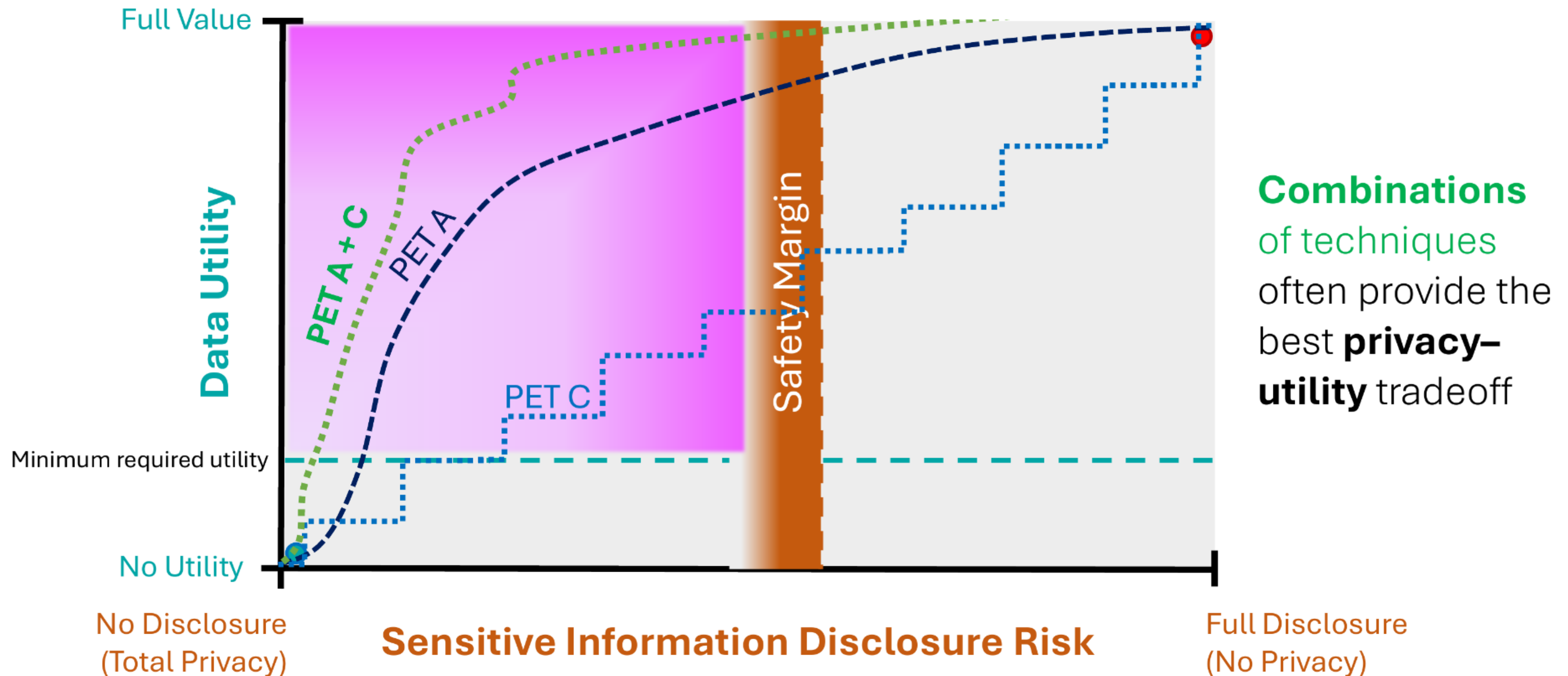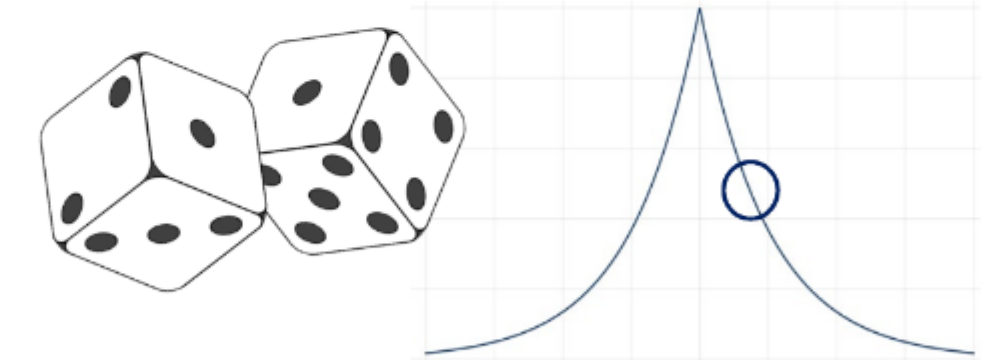Source: Adapted from Opening Report Figure 1: Privacy – utility tradeoff curve.

37

# Example: Combining PETs

| Query | Location |
|---|---|
| Mexican restaurant | **DC 20500** |
| resturant mexican | **DC 20500** |
| Mexican resturant | **DC 20500** |
| mexican history | DC 20500 |

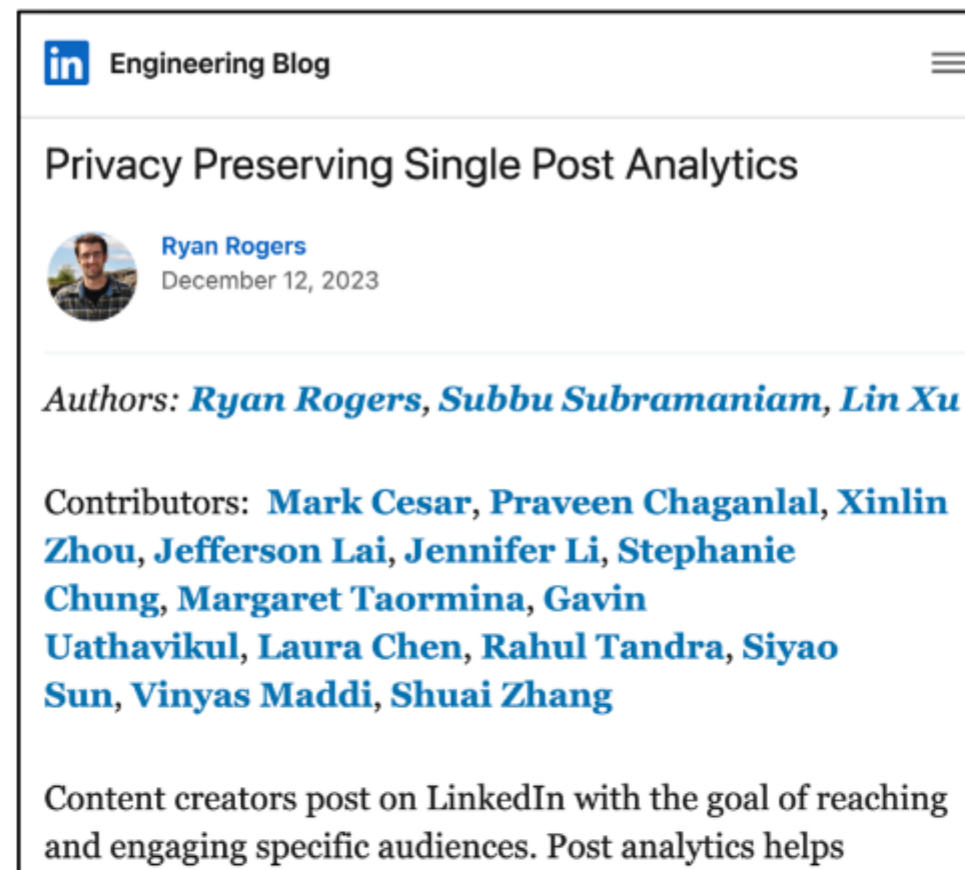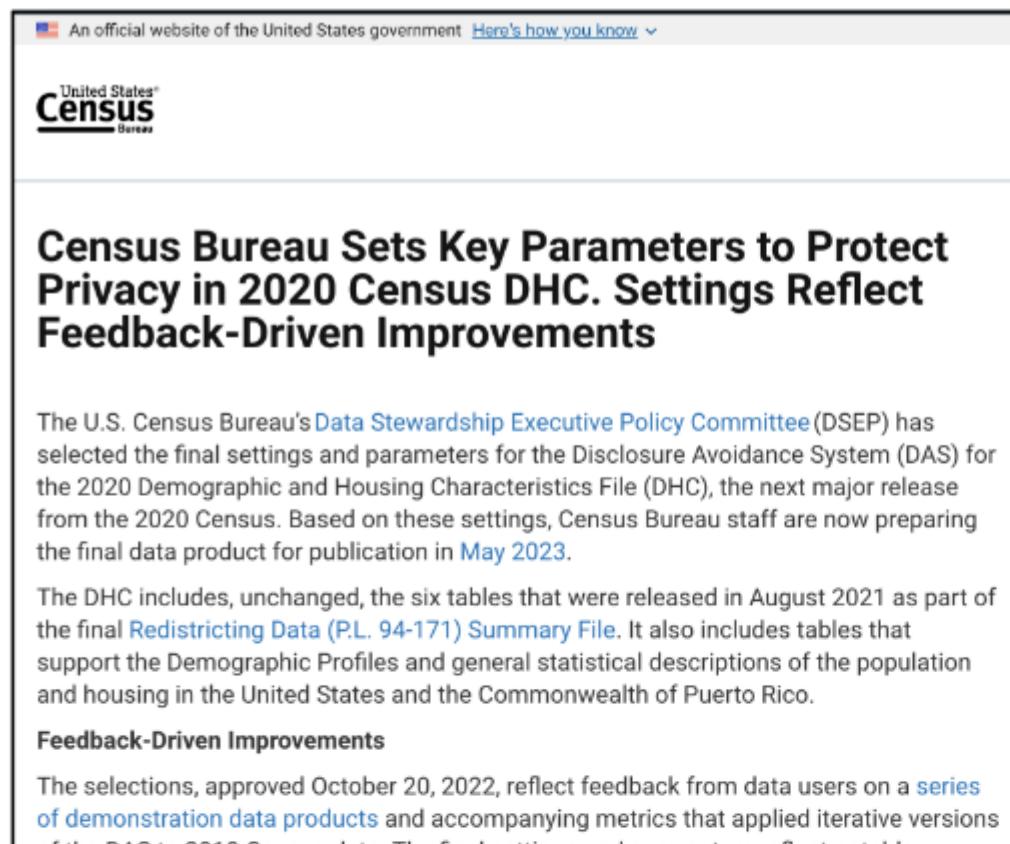**Generalization** to select (query, location)

**Differential Privacy Noise** to release statistics

| Query | Location | User Behaviors | |
|---|---|---|---|
| Mexican restaurant | DC 20500 | Count | 631.52 |
| | | Clicks | 472.24 |
| | | Average Time (s) | 2.24 |
| | | Abandoned | 18.02 |
| | | districttaco.com | 83.24 |
| | | dlenadc.com | 45.29 |
| | | mividamexico.com | 21.20 |
| | | ... | |

# Many Organizations Balance Privacy and Utility

# Google Has Experience Balancing Privacy and Utility



**SeDS Engineering Working Group**

Confidential

Privileged and Confidential

**DP for SeDS**

created: Jul 20, 2022
last updated: Jul 20, 2022
author: Dennis Kraft, Alex Kulesza, Sergei Vassilvitskii, Rachel Wei, Matthew Jagielski
status: WIP

"Over the years, we have **gained valuable experience** with DP, how it **translates to privacy policy** and how to implement it technically. Moreover, we have developed a **mature set of tools** to deploy DP quickly and efficiently."

**Robust privacy guarantees:** DP allows us to make strict and principled statements about privacy. If we enforce a certain DP specification, is it mathematically impossible to extract more information from the data than intended. This is particularly important when sharing data externally (as is the case for SeDS) given that we have limited control over the data after it has been released. Common sources of privacy issues DP is robustly protects against include:

**Internal Google Document**

**Differentially Private Stream Processing at Scale**

Bing Zhang[1], Vadym Doroshenko[†1], Peter Kairouz[†3], Thomas Steinke[†2], Abhradeep Thakurta[†2], Ziyin Ma[1], Eidan Cohen[1], Himani Apte[1], Jodi Spacek[1]
[1]Google
[2]Google DeepMind
[3]Google Research
{zhangbing,dvadym,kairouz,steinke,athakurta,ziyinma,eidanch,himaniapte,jodes}@google.com

**ABSTRACT**

We design, to the best of our knowledge, the first differentially private (DP) stream aggregation processing system at scale. Our system – *Differential Privacy SQL Pipelines (DP-SQLP)* – is built using a streaming framework similar to Spark streaming, and is built on top of the Spanner database and the F1 query engine from Google.

Towards designing DP-SQLP we make both algorithmic and systemic advances, namely, we (i) design a novel (user-level

called Differential Privacy SQL Pipelines (DP-SQLP), and make algorithmic advances along the way to cater to the scalability needs of it. DP-SQLP is implemented using a streaming framework similar to Spark streaming [52], and is built on top of the Spanner database [12] and F1 query engine [43] from Google. We also present production applications with two use cases in Section 6. The first is a real world use case that deploys DP-SQLP in *Google Shopping* to generate streaming page-view counts. The second applies the streaming DP

"In terms of **data utility** after adopting DP-SQLP, we were able to retain 59% of the page-view.... to **99.9% for pages with an average view rate of 60 views/hour**. When comparing noised impression counts with the raw counts, the **relative error** is around 11%.... to ensure **user level DP guarantee**, per day. We use $\varepsilon = 1$ for ...."

**Google Research Paper**

**REDACTED FOR PUBLIC FILING**

# Opinion 3: Data at Issue Can Be Shared Safely

**1** There are well-established **privacy-enhancing techniques** that can be used to protect sensitive information.
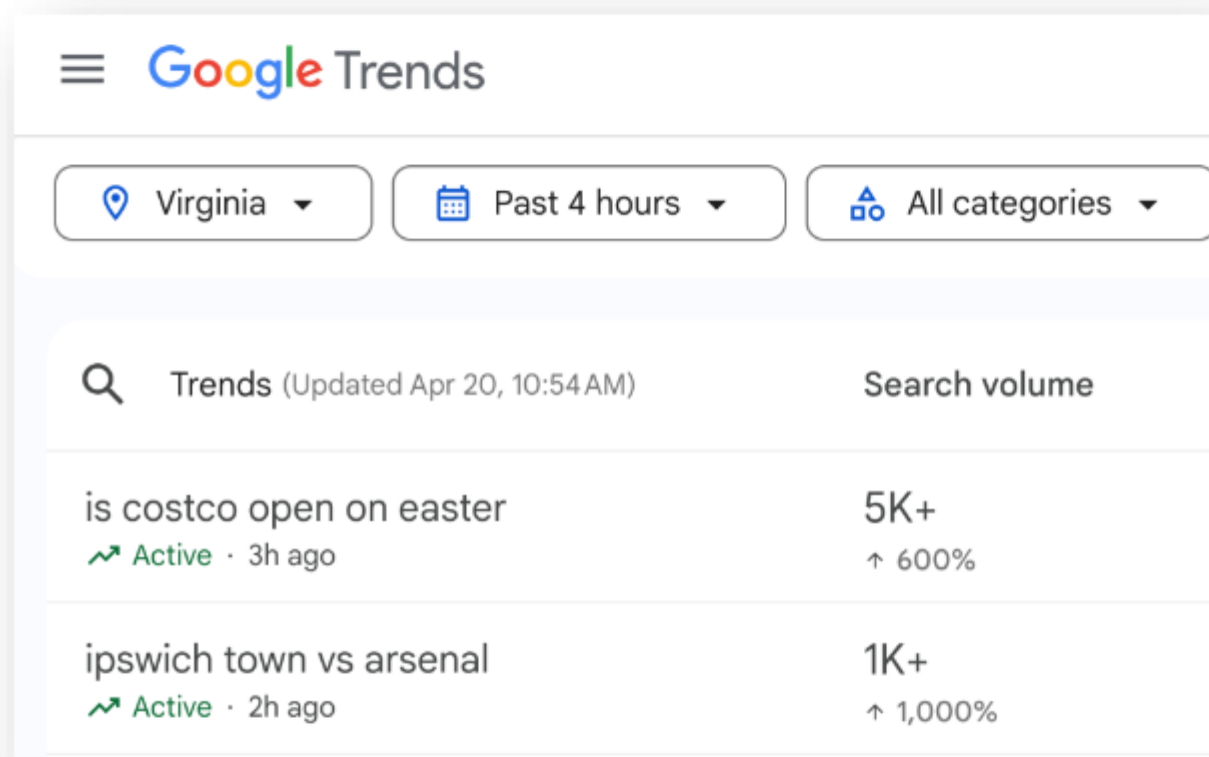
**2** Many organizations, including Google, **safely release sensitive data** by using privacy-enhancing techniques.

**3** Google can share the data at issue in a way that **assures privacy** while **providing utility**.

**REDACTED FOR PUBLIC FILING**

# Google Currently Uses PETs to Release Similar Data
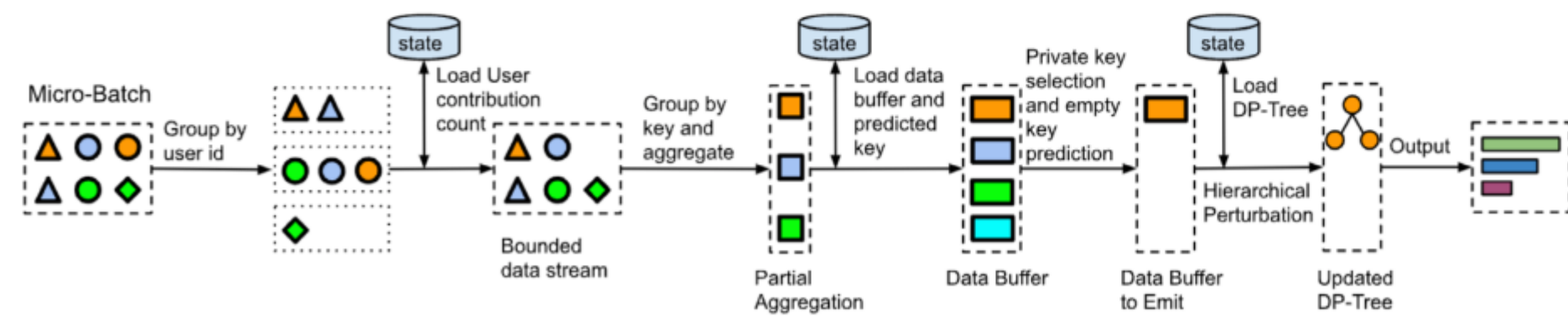
## Search Queries



**Google Trends**

Covid Symptoms, Vaccination Insights, …

## User Interactions



**Google Shopping**

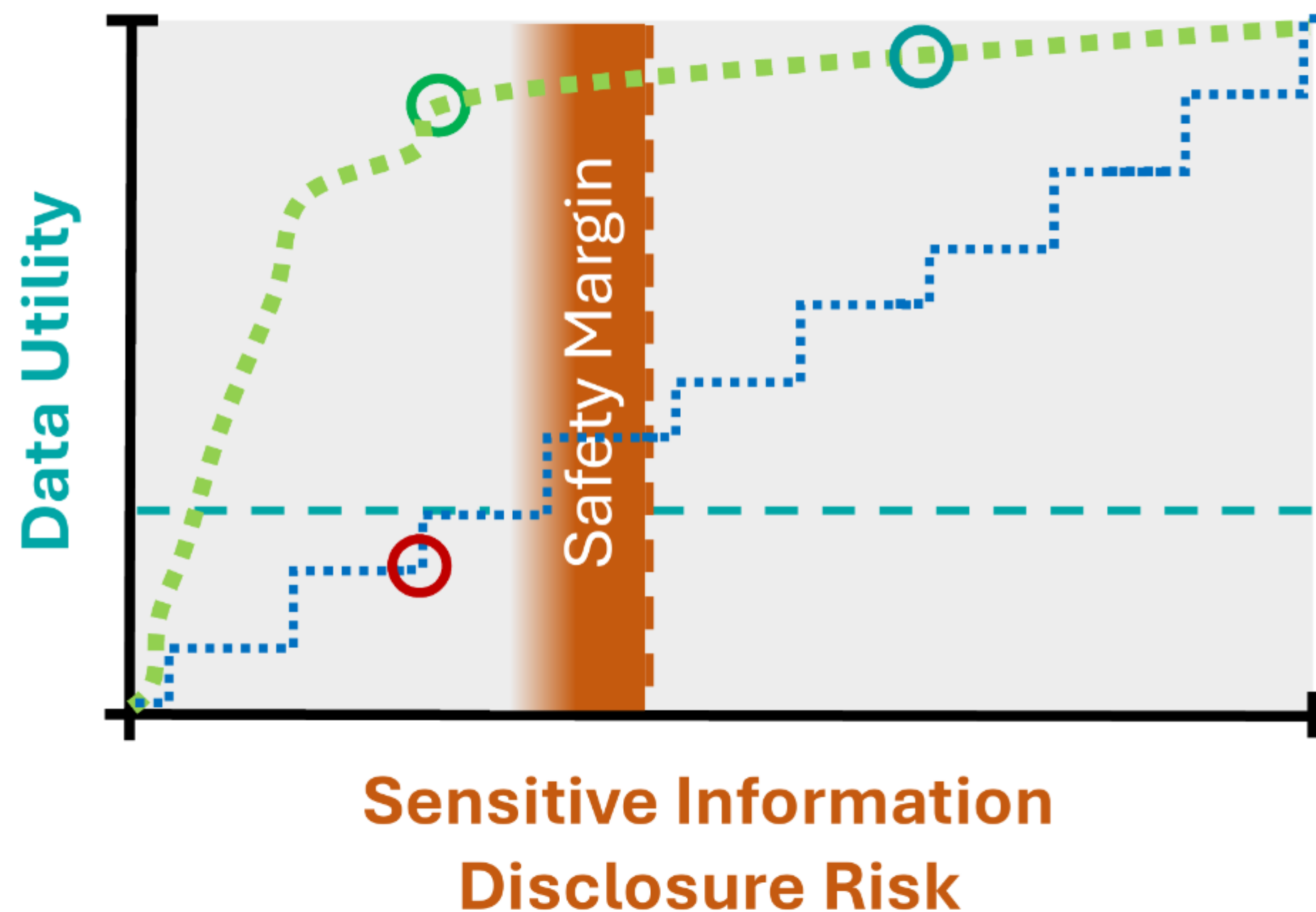## Advertising Data

Private-Set Intersection, Analytics

## Real-time

Google Trends, Google Shopping

## Enormous Scale

**Plume** (Trillions of records with DP)

# Implementing the Data Sharing Remedy



**Data Utility** / **Safety Margin** / **Sensitive Information Disclosure Risk**

The **Technical Committee** with understanding of **intended uses** and **data content** can assess use of privacy-enhancing techniques and parameters for an appropriate **privacy–utility** tradeoff.

# Google's Expert Agrees: Data Can Be Shared Safely

## Dr. Culnane's Deposition

Q.      Dr. Culnane, you believe that it is possible for Google to share what you call the DOJ search data by applying privacy-enhancing techniques to achieve suitable privacy safeguards, don't you?

A.      Yes.

Q.      Do you have any opinion as to whether it is technologically feasible to share the DOJ search data as Plaintiffs describe in Plaintiffs' Proposed Final Judgment?

A.      The subject of my report is looking at the ability to do that safely, so there is an opinion as -- if it is correctly protected, and in my view, if you protect personal data as opposed to PII, then you can anonymize the dataset. If you successfully do that, then you can protect privacy by doing that, yes.

# Conclusion

**1** There are well-established **privacy-enhancing techniques** that can be used to protect sensitive information.

**2** Many organizations, including Google, **safely release sensitive data** by using privacy-enhancing techniques.

**3** Google can share the data at issue in a way that **assures privacy while providing utility.**