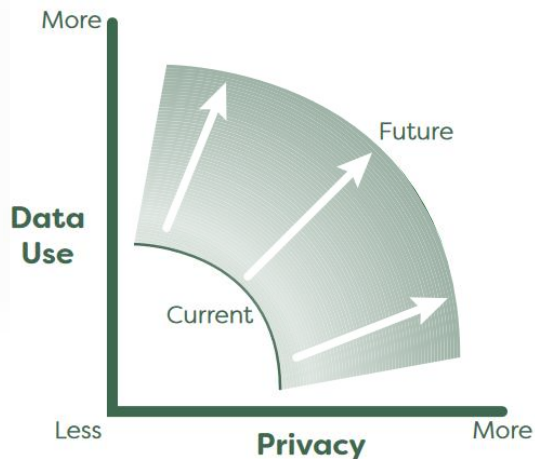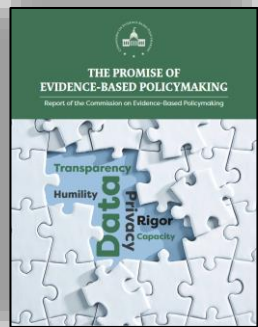# CHALLENGES TO PRIVACY AND CONFIDENTIALITY ACROSS THE BLENDED DATA LIFECYCLE

Mayank Varia

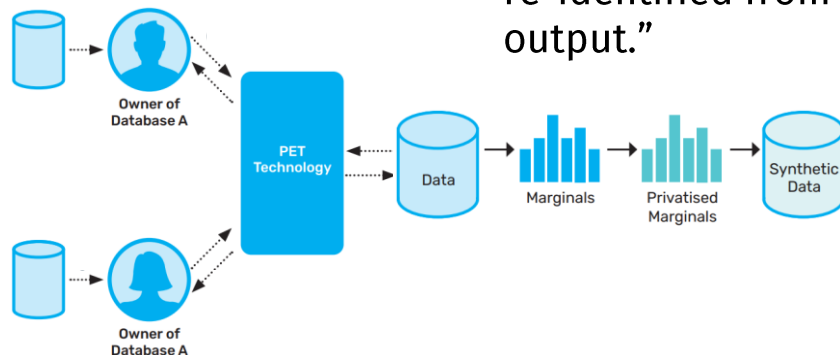Boston University

May 22, 2023

# WHY USE PRIVACY TECH?



"Privacy, when implemented correctly, fosters more information sharing, not less."

– Marc Groman

"The benefits of using data for official statistics can be realized while minimizing privacy risks to those entrusting sensitive data to National Statistics Offices."

# INPUT PRIVACY

# OUTPUT PRIVACY

"Allow two or more parties to submit data into a calculation without the other respective parties seeing data in clear."

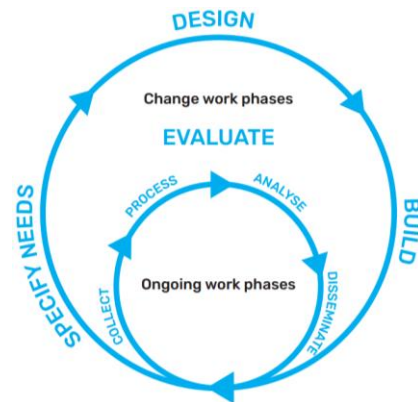"Generally known as statistical disclosure control, [it] aims to conceal sensitive individual data from being identified or re-identified from the disseminated output."
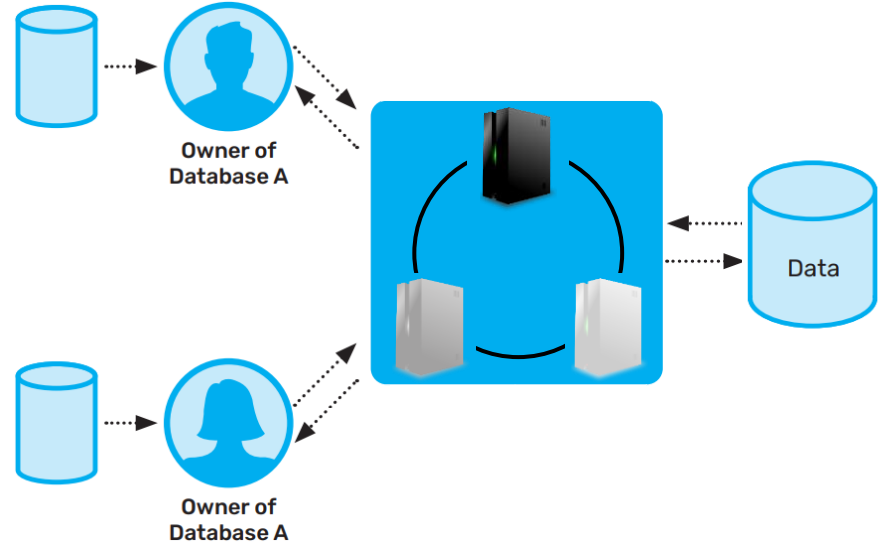


Selected examples:

- Secure multi-party computation
- Homomorphic encryption
- Zero knowledge proofs

Examples:

- Differential privacy
- Synthetic data generation

Image source: UN PET Guide

# SECURE MULTI-PARTY COMPUTATION (SMC)

"Enables different participating entities in possession of private sets of data to link and aggregate their data sets … **without transferring or otherwise revealing any private data** to each other or anyone else."
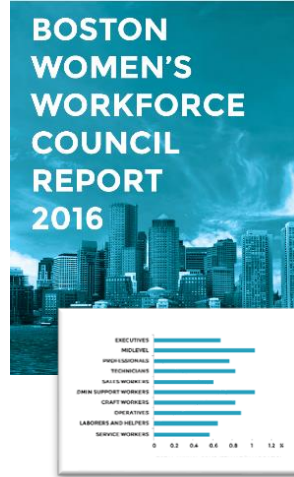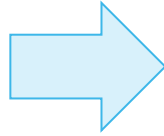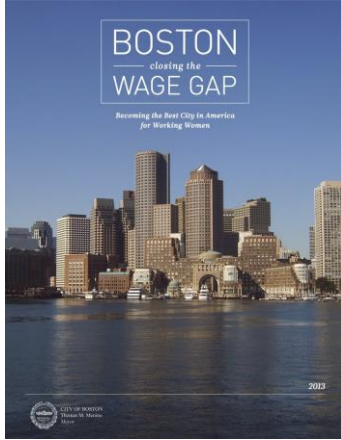
– 2022 U.S. Senate bill S.3952

# SMC DEPLOYMENTS IN THIS TALK

1. Boston Women's Workforce Council (non-blended data)

2. Massachusetts child advocacy study (blended data)

# BOSTON WAGE GAP STUDY
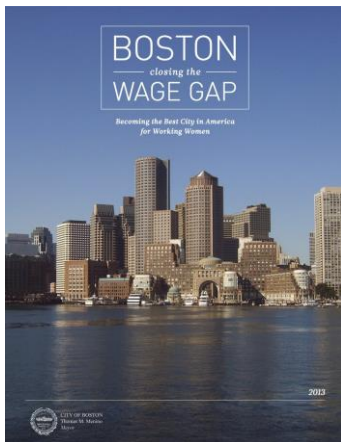


Lesson: Deploying SMC is a policy decision, promote from the top down

# DATA CONTRIBUTORS



"Employers agree to anonymously contribute their wage data with the Council, which then creates a snapshot of what the wage gap looks like in our city."



Lesson: Easier to deploy SMC when all data contributors have similar interest in the output, and are similarly concerned about input privacy
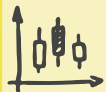
# Boston Women's Workforce Council

100% Talent Data Submission
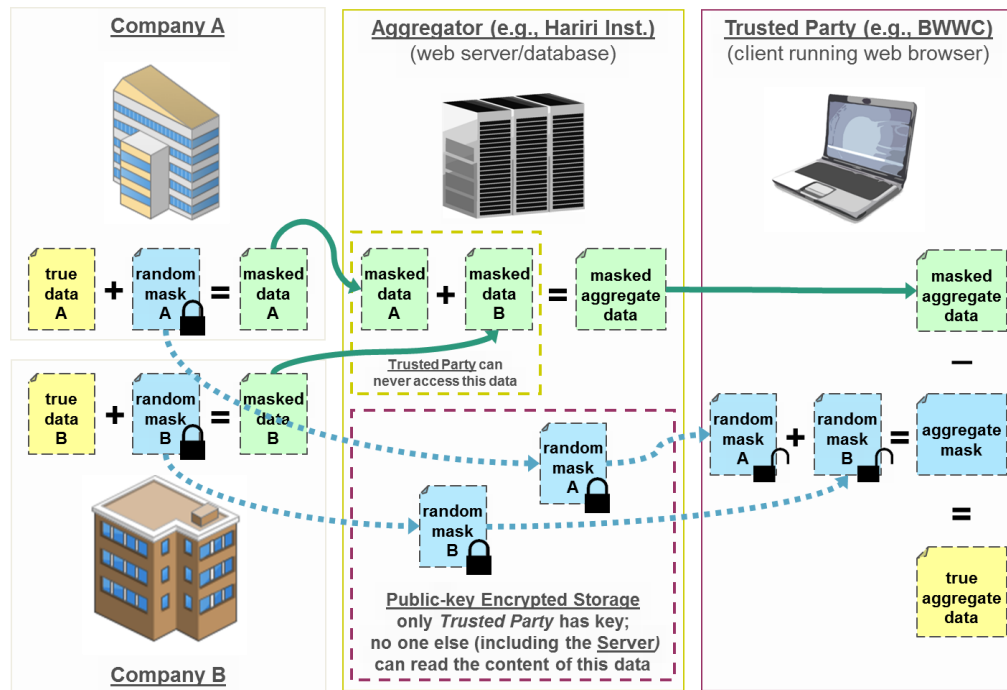
## Number Of Employees

| | Hispanic or Latinx | | White | | Black/African American | | Native Hawaiian or Pacific Islander | | Asian | | American Indian/Alaska Native | | Two or More Races (Not Hispanic or Latinx) | | Unreported | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Female | Male | Female | Male | Female | Male | Female | Male | Female | Male | Female | Male | Female | Male | Female | Male |
| Executive/Senior Level Officials and Managers | | | | | | | | | | | | | | | | |
| First/Mid-Level Officials and Managers | | | | | | | | | | | | | | | | |
| Professionals | | | | | | | | | | | | | | | | |
| Technicians | | | | | | | | | | | | | | | | |
| Sales Workers | | | | | | | | | | | | | | | | |
| Administrative Support Workers | | | | | | | | | | | | | | | | |

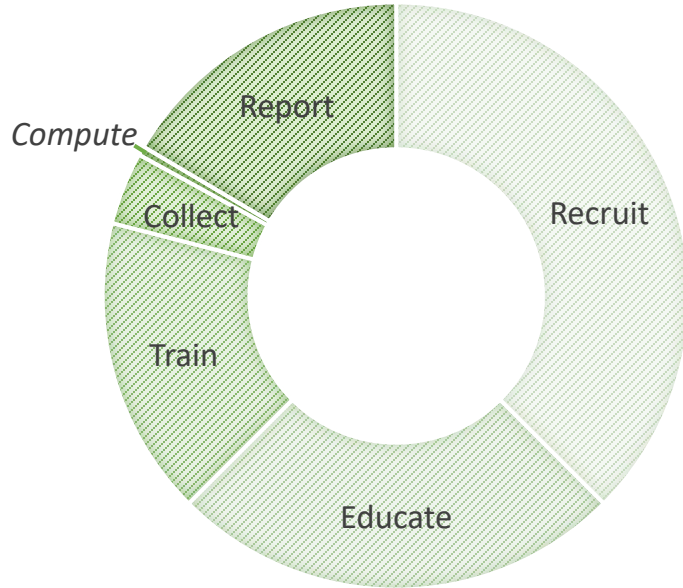Lesson: Use standardized data formats & easily accessible applications
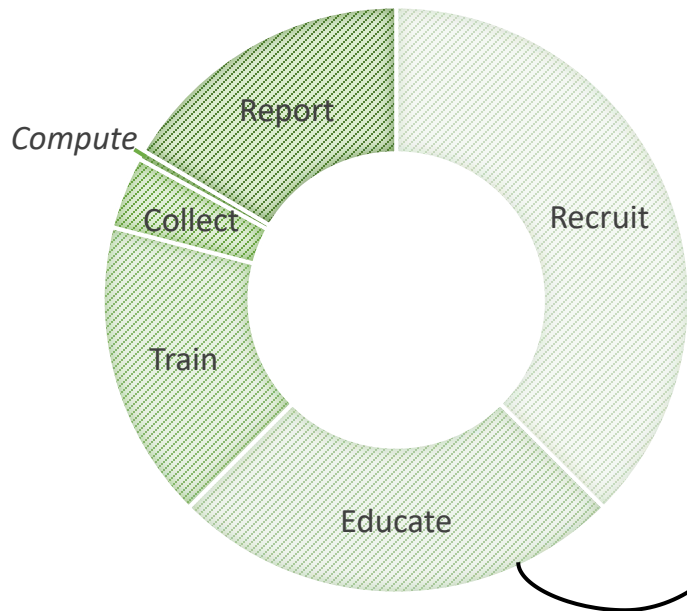
# SIMPLICITY & TRANSPARENCY DRIVE ADOPTION



👓 Lesson: SMC is simple to explain, and doing so improves trust + adoption

# TIMELINE



Lesson: Time is spent on finding and helping people to work with data, not computer science metrics like computational efficiency

# EDUCATION & TRAINING



Report

Compute

Collect

Recruit

Train

Educate



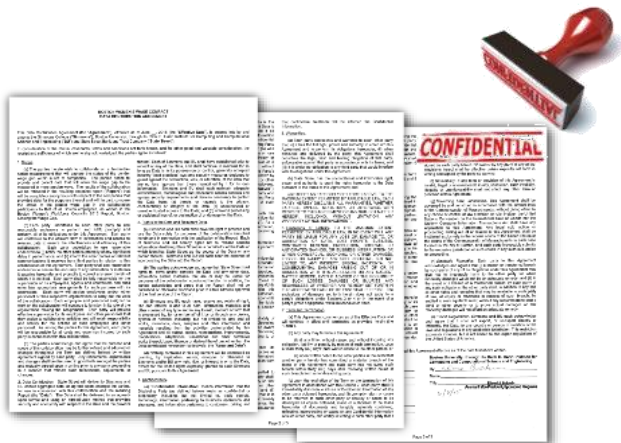**WomensWrkfrceCouncil**
@BostonWomenWork

Follow

#100PercentTalent Compact Signers learning about revolutionary MPC technology over lunch. How's that for a working lunch?

Lesson: Identify key people to convince, together with domain experts

# NDAs & LEGAL CONSIDERATIONS

**Definitional question**
Do encodings count
as personal data?

**Process question**
Does computing
constitute disclosure?

**Liability question**
Who should be
blamed for an error?

Lesson: Consider the legal ramifications of data disclosure & processing
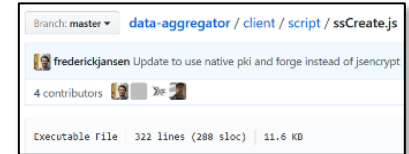
# VALUE PROPOSITION

**BWWC**
Data quality

**HR Personnel**
Accessibility

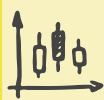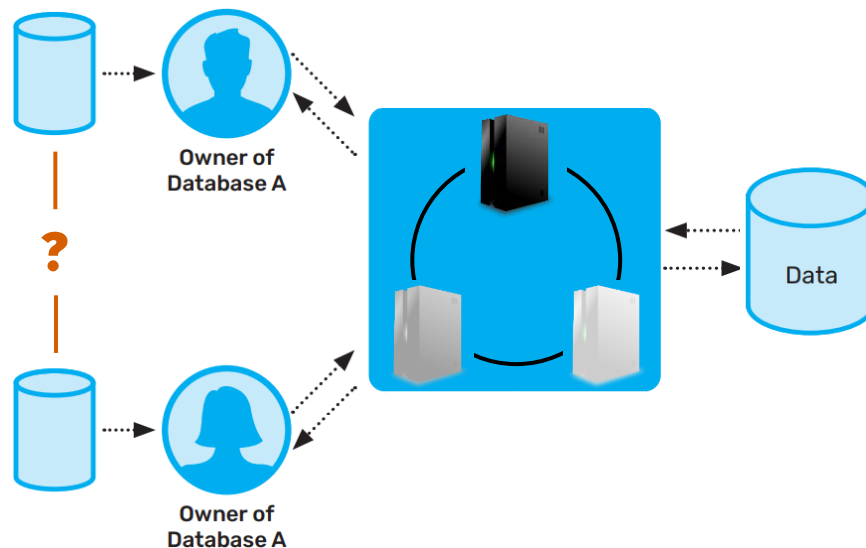**Lawyers**
Liability

**IT Personnel**
Comprehension



Lesson: Provide value to all stakeholders involved in the data analysis

# CHILD ADVOCACY STUDY

Goal: make data-driven recommendations for reforming the juvenile justice system

Blend datasets containing data on:

- Vulnerable groups
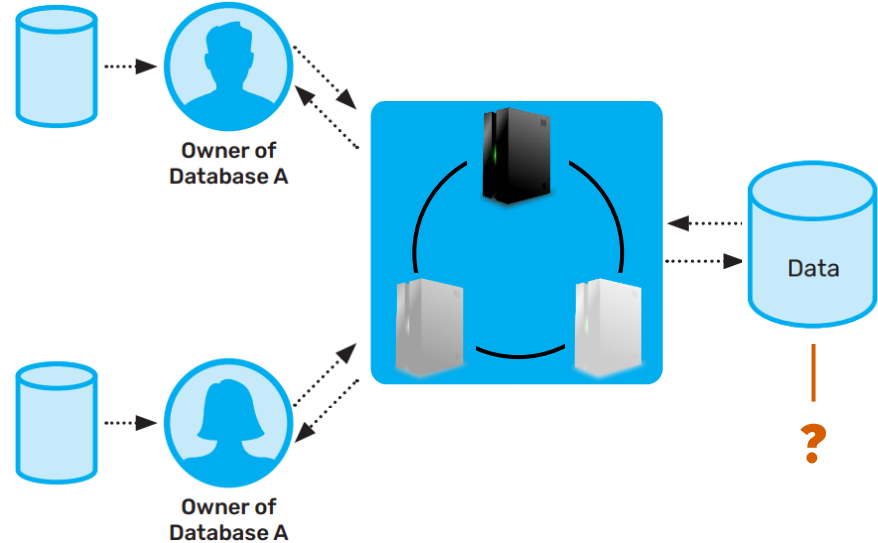- Educational success
- Judicial outcomes
- ...



Lesson: Challenging to conduct any data analysis (with or without input privacy) in the absence of standardized data formats and common fields

# CHILD ADVOCACY STUDY

Goal: make data-driven recommendations for reforming the juvenile justice system

Blend datasets containing data on:

- Vulnerable groups
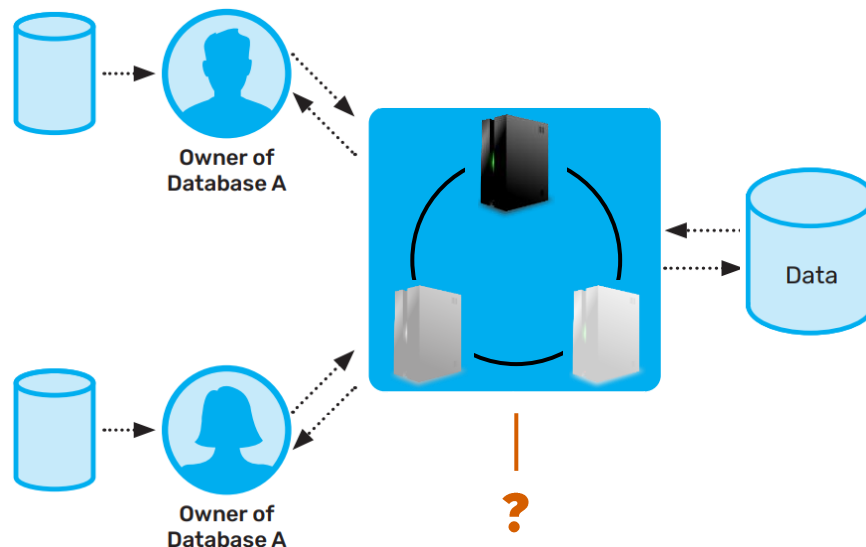- Educational success
- Judicial outcomes
- ...



Lesson: Challenging to build consensus for data analyses when each potential contributor has its own policy goals and interest in the output

# CHILD ADVOCACY STUDY

Goal: make data-driven recommendations for reforming the juvenile justice system

Blend datasets containing data on:
- Vulnerable groups
- Educational success
- Judicial outcomes
- …



Owner of Database A

Owner of Database A

Data

?

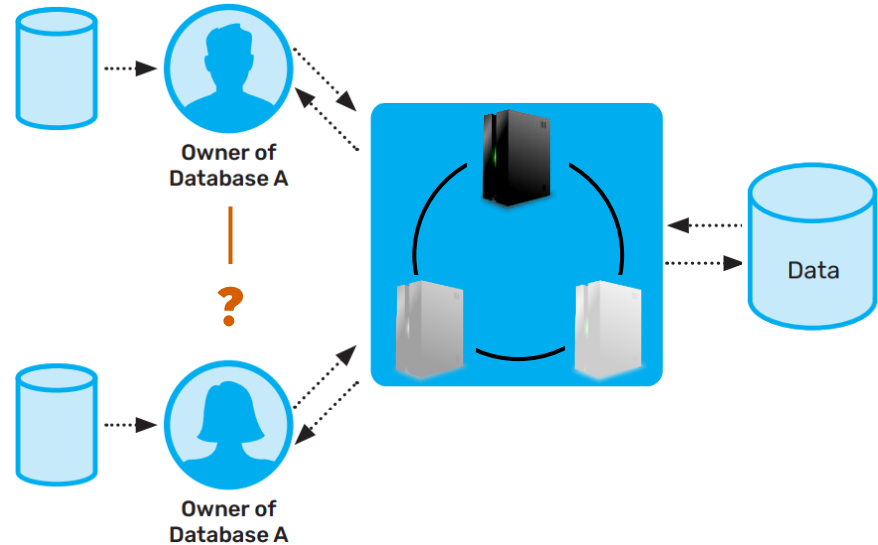Lesson: Sometimes privacy questions are masking deeper concerns

# CHILD ADVOCACY STUDY

Goal: make data-driven recommendations for reforming the juvenile justice system

Blend datasets containing data on:

- Vulnerable groups
- Educational success
- Judicial outcomes
- …



Owner of Database A

Owner of Database A

Data

?

Lesson: Ultimately, can only move forward with a data analysis if the result provides value to all stakeholders involved
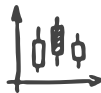
# SUMMARY: LESSONS LEARNED

| WHO | WHAT | HOW |
|---|---|---|
| Promote from the top down | Use standardized data, schemas, apps | Provide value to all stakeholders |
| Find orgs with interest in the result | Identify key people to convince | Simple, transparent tech drives adoption |
| People, not CPU, dominate runtime | Consider legal ramifications | Improve both data privacy and utility |