



# Demystifying Privacy Preserving/Enhancing Technologies

---

**Stephanie Straus**

Massive Data Institute  
McCourt School of Public Policy  
Georgetown University



## What We Do

- **Convene** and advocate on behalf of communities that are sharing and using cross-sector data for good
- **Connect** to innovations, best practices, and research and funding opportunities that support ethical data sharing
- **Consult** with data sharing collaborations to build the human and technical capacity to share data and improve lives

## Why We Do It

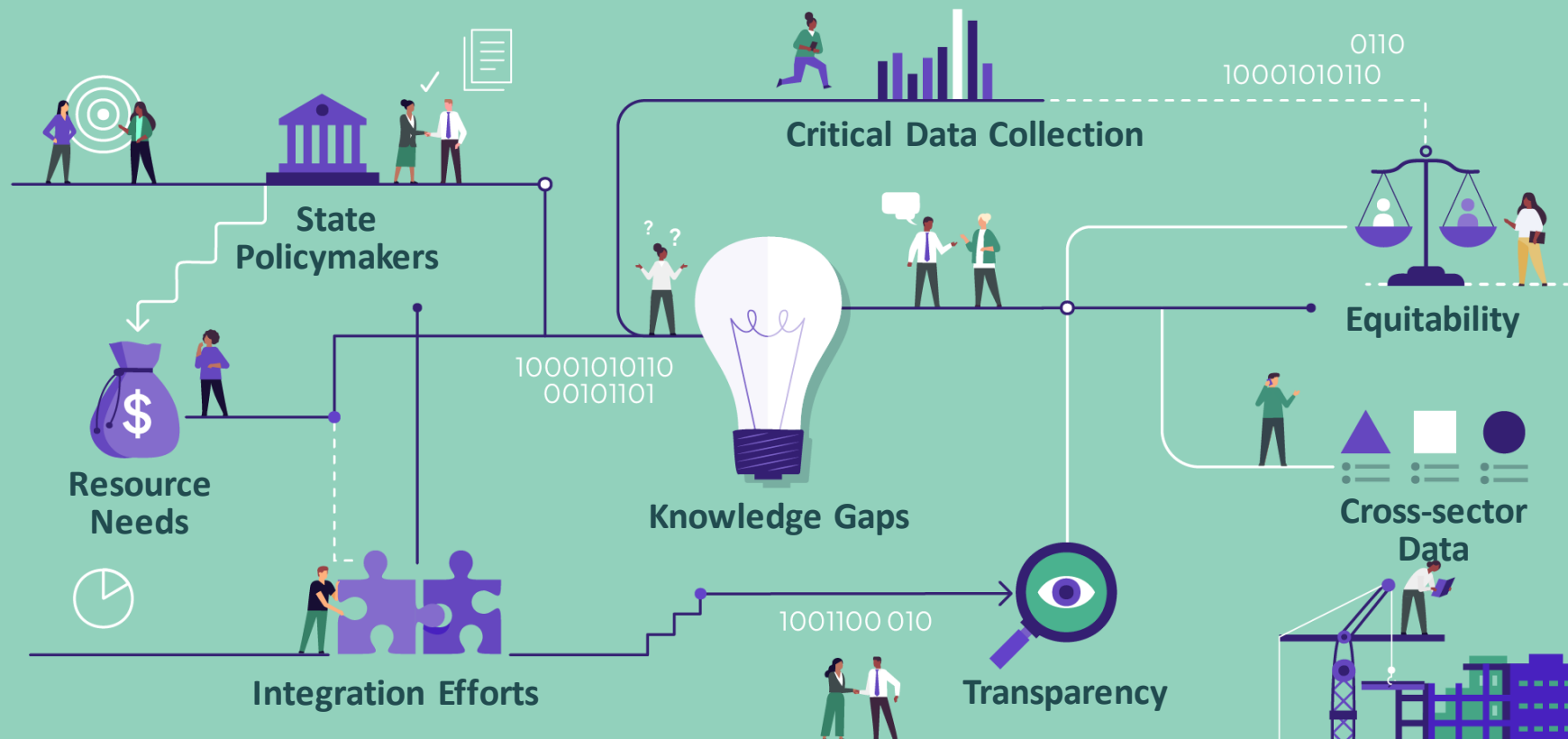
When communities bring together cross-sector data safely and responsibly, policy-makers, practitioners, and schools are better equipped to:

- Understand the complex needs of individuals and families
- Allocate resources where they're needed most to improve services
- Measure long-term and two-generation impacts of policies and programs
- Engage in transparent, shared decision-making about how data should (and should not) be used

[www.aisp.upenn.edu](http://www.aisp.upenn.edu)



The Data Integration Support Center (DISC) at WestEd provides expert integrated data system planning and user-centered design, policy, privacy, and legal assistance for public agencies nationwide.



# Massive Data Institute (MDI)



The Massive Data Institute (MDI) at Georgetown's McCourt School of Public Policy focuses on the secure and responsible use of data to answer public policy questions.

MDI works with researchers in government, academia, and industry to solve societal-scale problems using novel and traditional large-scale data sources.

MDI's strategic partnerships promote community and innovation across the health, social, computer, and data sciences.

# Our roles



## We are:

Data evangelists

Connectors, community builders,  
thought partners, cheerleaders,  
and data sharing therapists

Focused on ethical data use  
for policy change



## We are not:

Data holders or intermediaries

A vendor or vendor recommenders

Focused on academic research

# LEGAL DISCLAIMER



- Not Legal Advice
- Training will only cover **federal law**
- Laws change, this is based on the law at the time of the training
- Consult your general counsel for specific legal questions

# Essential Questions



What are privacy-enhancing (PET) or preserving technologies (PPT)?



What are the benefits and risks of these technologies and tools?



What factors or questions should be considered before building, procuring, or using PETs?

# Modern Data System Technical Capabilities



## Link and Resolve

Resolve and maintain unique master IDs for each mastered entity

Link different data sets through predefined or discovered identifiers



## Deliver Data Sets

Generate and transfer individual-level data sets to requestors securely



## Deliver Visualizations and Analyses

Deliver analysis to the public without requiring logins, through interactive visualizations and summary data

Deliver sensitive analysis results securely to an authorized organization

Deliver sensitive analysis results securely to authorized individuals



## Analyze, Interpret, and Predict

Statistical models are developed and applied to derive meaning from data

Data are used for predictions using statistical and machine learning techniques



# Current Privacy Protection Methods



## Limit Data Access

Severely restrict or eliminate access to the data



## Legal and Contractual Safeguards

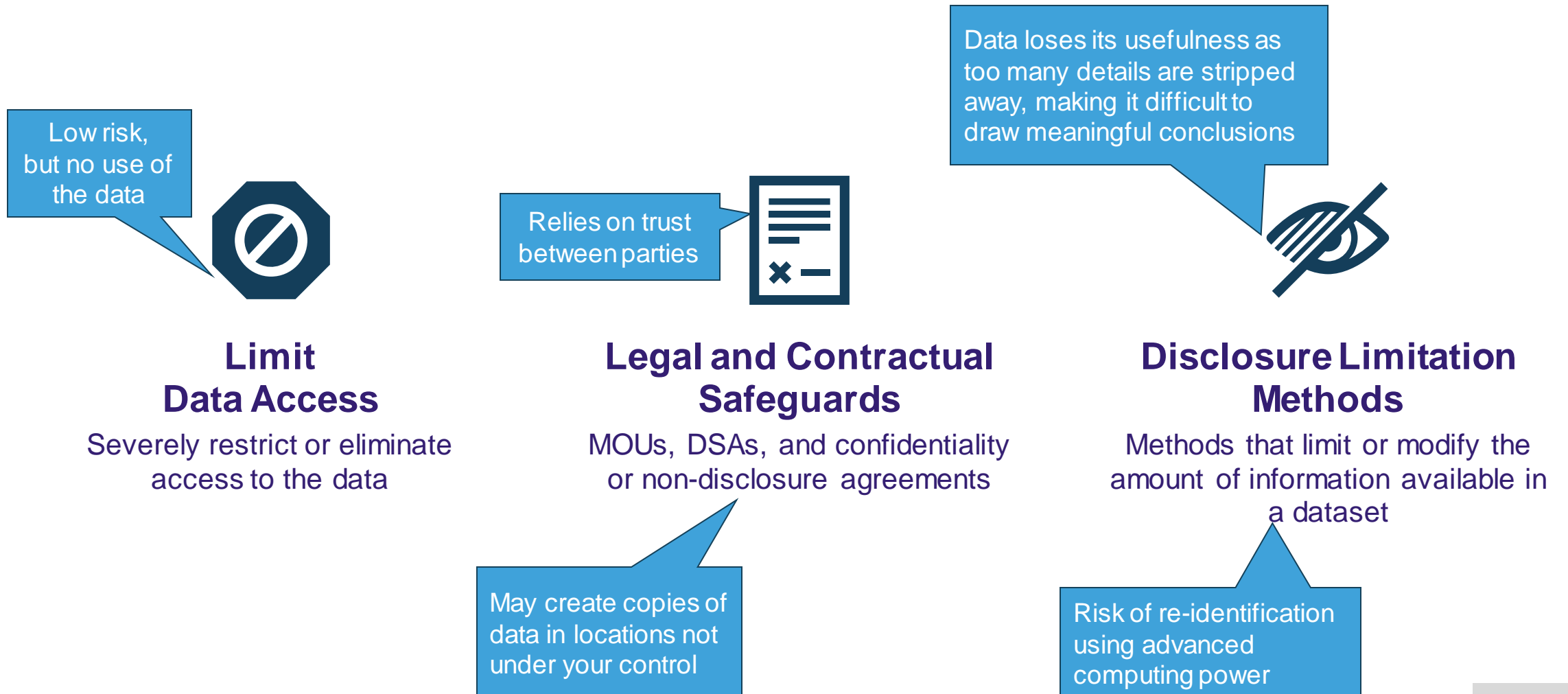
MOUs, DSAs, and confidentiality or non-disclosure agreements



## Disclosure Limitation Methods

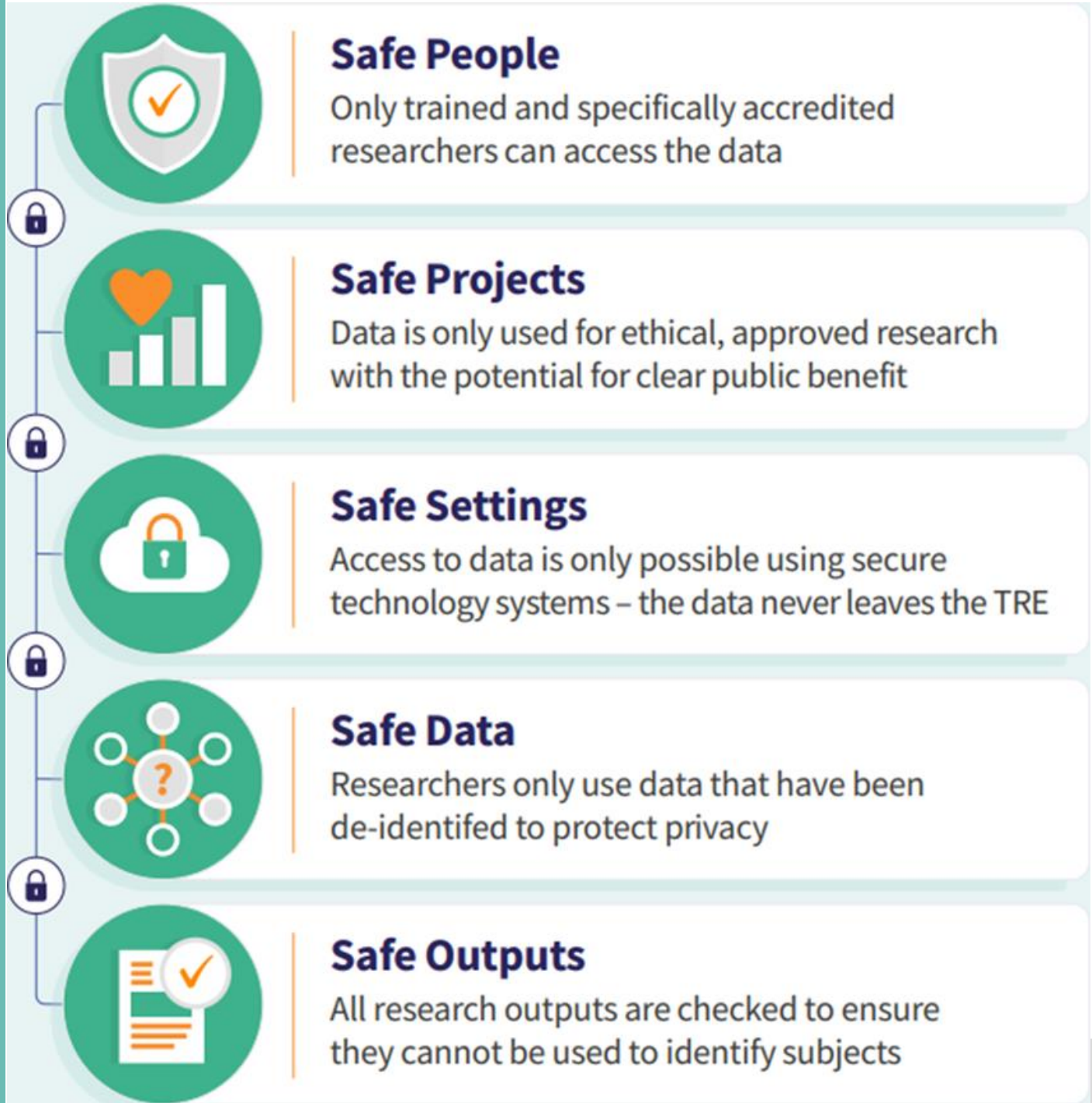
Methods that limit or modify the amount of information available in a dataset

# Current Privacy Protection Methods



# The Five Safes Framework

Scientific Figure on ResearchGate from [Towards a standardised cross-sectoral data access agreement template for research: a core set of principles for data access within trusted research environments](#)  
[CC BY 4.0]



# What are PETs?



# What are Privacy Enhancing Technologies (PETs)?



**PETs are safer and more secure ways to analyze, link, and share data**



Cryptographic techniques that increase data protection while allowing for greater data utility



Also known as Privacy Preserving Technologies (PPT)



Can enhance how data are analyzed and/or published

# Five key components of successful integrated data systems



## Governance

The people, policies, and procedures that support how data are used and protected.



## Legal

The legal framework supports the purpose for data sharing, documents the legal authority of the host organization, and ensures that data sharing complies with all federal and state statutes.



## Technical

Technical components are created to support the core purpose



## Capacity

Staff, relationships, and resources that enable an effort to operate governance, establish legal authority, build technical infrastructure, and demonstrate impact.



## Impact

All components of quality – governance, legal agreements, technical tools, staff capacity – exist to drive impact.

**PETs**

# Five Safes + PETs



Distrusting parties.  
Less 'eyes' on the sensitive info

Set the queries, less inappropriate use or unauthorized access

Access controls, encryption in location, less data leaving 'home'

Obfuscate or remove PII, encrypt entire dataset

'Noise up' the data, synthesize;  
reduce re-ID risk

# How PETs Address Data Governance Issues



## Input Privacy

- Secure hashing
- Trusted execution environments or secure enclaves
- Intermediaries
- Secure multiparty computation
- Homomorphic encryption

## Output Privacy

- Traditional SDL
- Differential privacy
- Private query server
- Synthetic data





# Secure Multiparty Computation (MPC)

**Secure multiparty computation (MPC):** the process by which two distrusting parties jointly compute a research query on their datasets, without ever seeing the other's underlying data.

- ☒ use for aggregate statistics or individual level data
- ☒ use for overlap in datasets
- ☒ no third party needs to see data/parties communicate directly

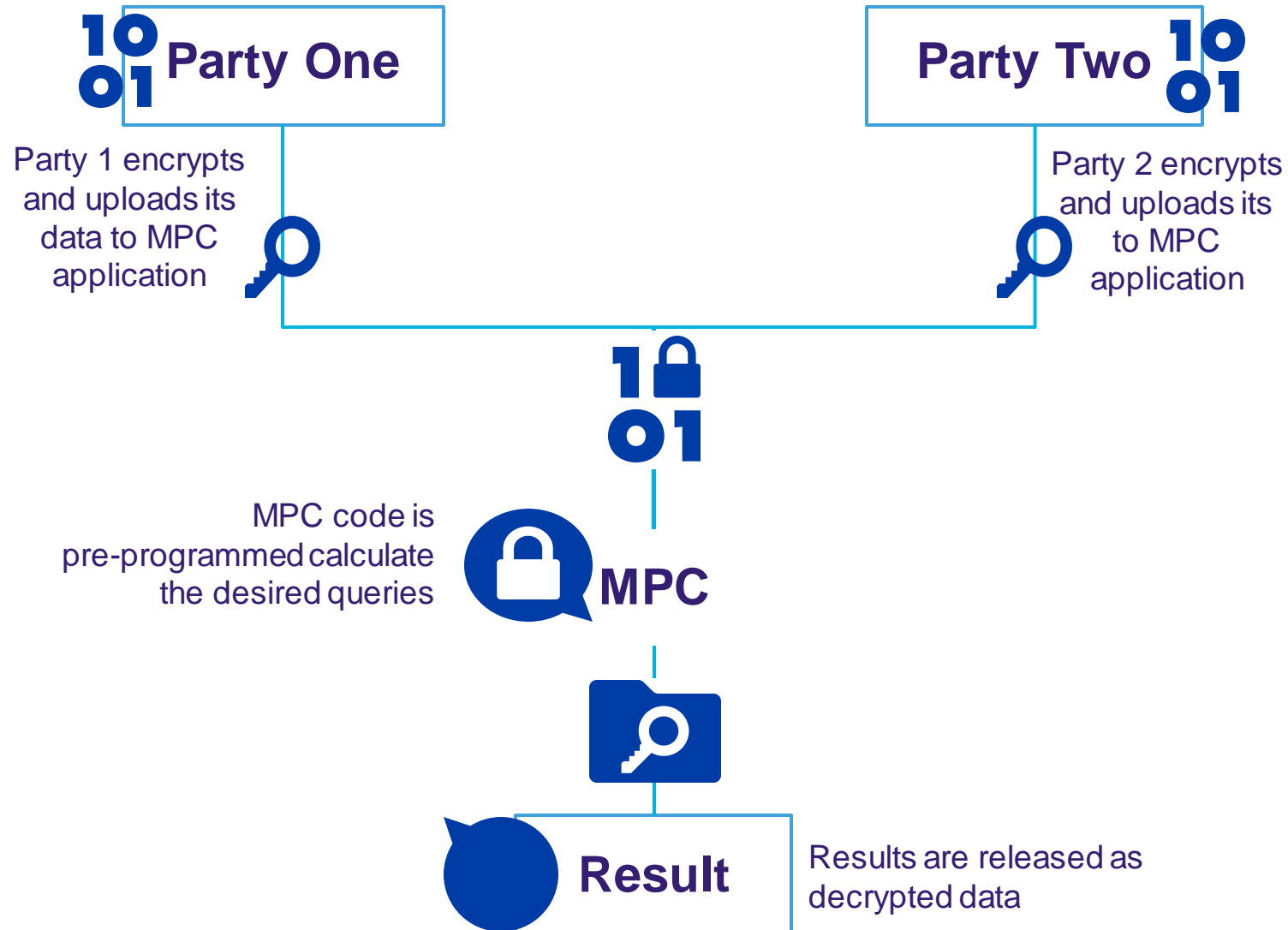
- ☐ time-consuming
- ☐ limited in operations/statistics
- ☐ requires careful data preparation
- ☐ does not address output privacy

## Examples:

- [Estonia](#)
- [Virginia](#)
- [our NCES demonstration](#)
- [Boston Women's Workforce](#)
- [Allegheny County Department of Human Services demonstration](#)
- DARPA and IARPA investments



# Secure Multiparty Computation (MPC)



# FAQ: Secure Multiparty Computation

Who can see the original data?

- Only the original data owners, before they upload their data to the MPC application.

Who has the decryption key(s)?

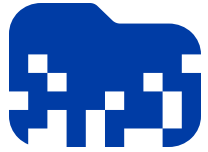
- Usually, the data owners, by design. Keys should never be shared.

What queries are being asked?

- Up to the data owners.

What would I see if intercepted messages?

- Ciphertext (gibberish).



# Secure Enclave

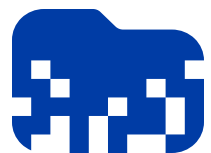
**Secure enclave:** a virtual computing workspace that enables authorized users to access sensitive data and securely conduct analysis.

- ☒ supports open-ended queries, not limited in operations/statistics
- ☒ external users cannot download or extract data without permission

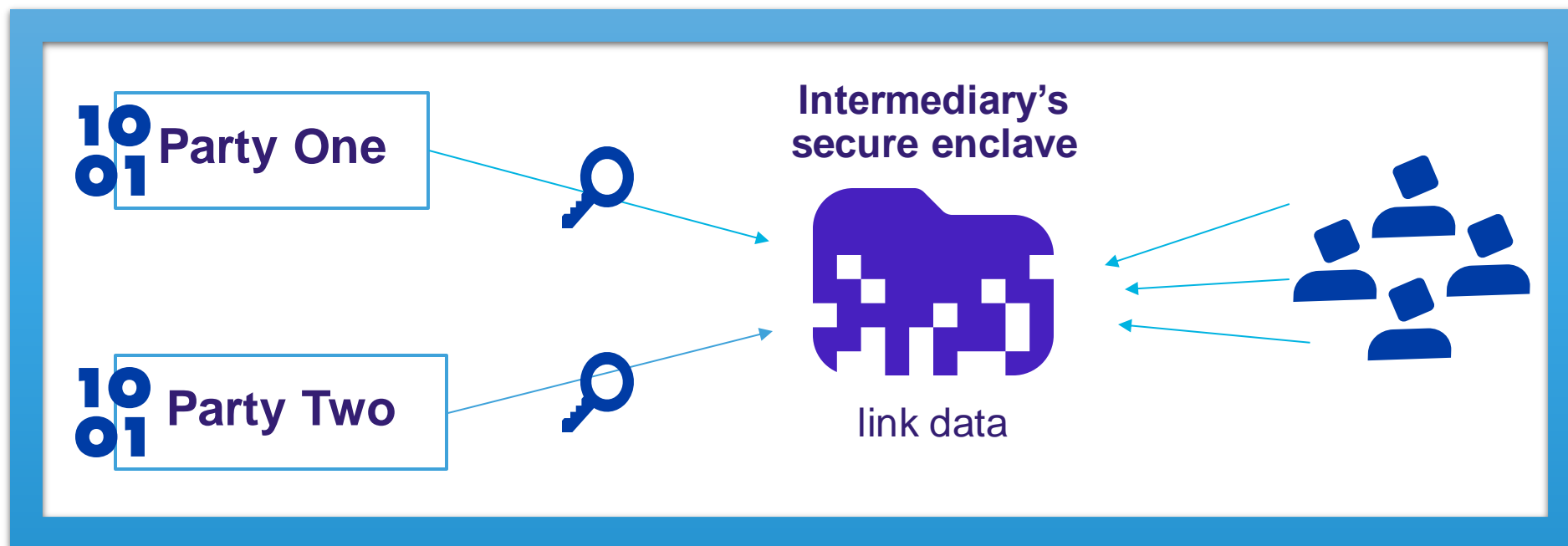
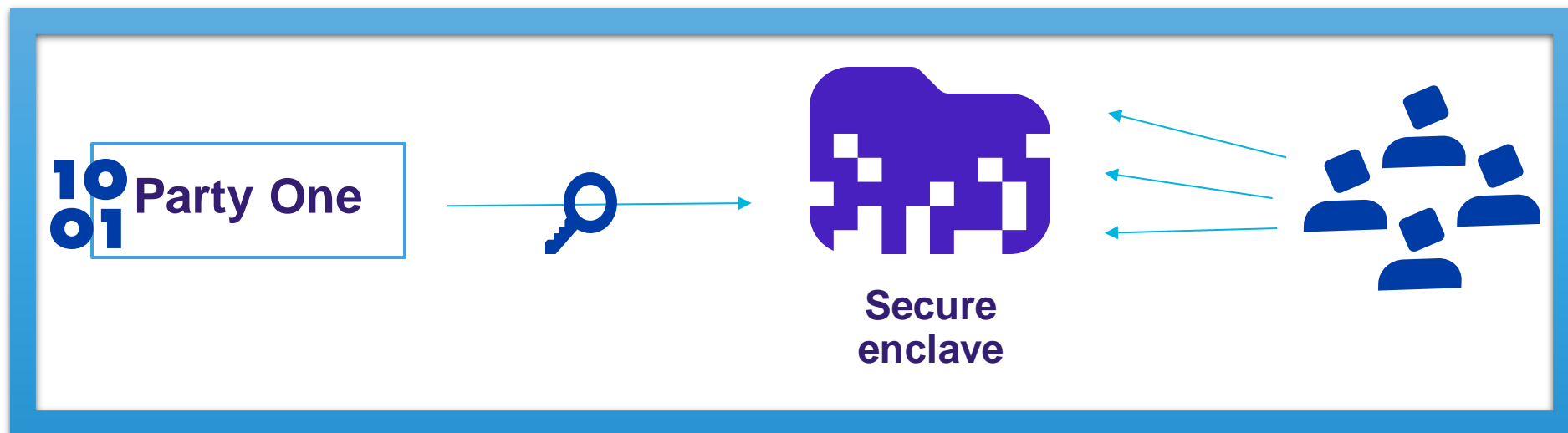
- ☐ trusted third party sees data (to de-identify, provision, etc.)
- ☐ external users see potentially sensitive data
- ☐ must layer output privacy protections

## Examples:

- [Washington State Education Research Data Center \(ERDC\)](#)
- [Coleridge Initiative](#)
- [Department of Justice](#) (internal)
- [USDA Economic Research Service](#)



# Secure Enclave



# FAQ: Secure Enclaves

What if I try to hack the enclave?

- Data are typically encrypted at rest. Moreover, TEEs store memory separate from computer's CPU.

What's a TEE?

- A Trusted Execution Environment (TEE) is a secure enclave, but with hardcoded security and access controls.

Secure Enclave	Trusted Execution Environment (TEE)
<ul style="list-style-type: none"><li>• manual security and access controls</li><li>• need 'eyes' on the data</li><li>• relies on authorized users following governance protocols</li></ul>	<ul style="list-style-type: none"><li>• automated security and access controls</li><li>• can manage the TEE 'blind,' without seeing the code or data</li><li>• relies on cryptographic verification</li></ul>

# #Secure Hashing

**Secure hashing:** an algorithm that replaces sensitive inputs with a random string of characters (hash) unique to each original record in the data. Can be used to protect passwords or link across datasets.

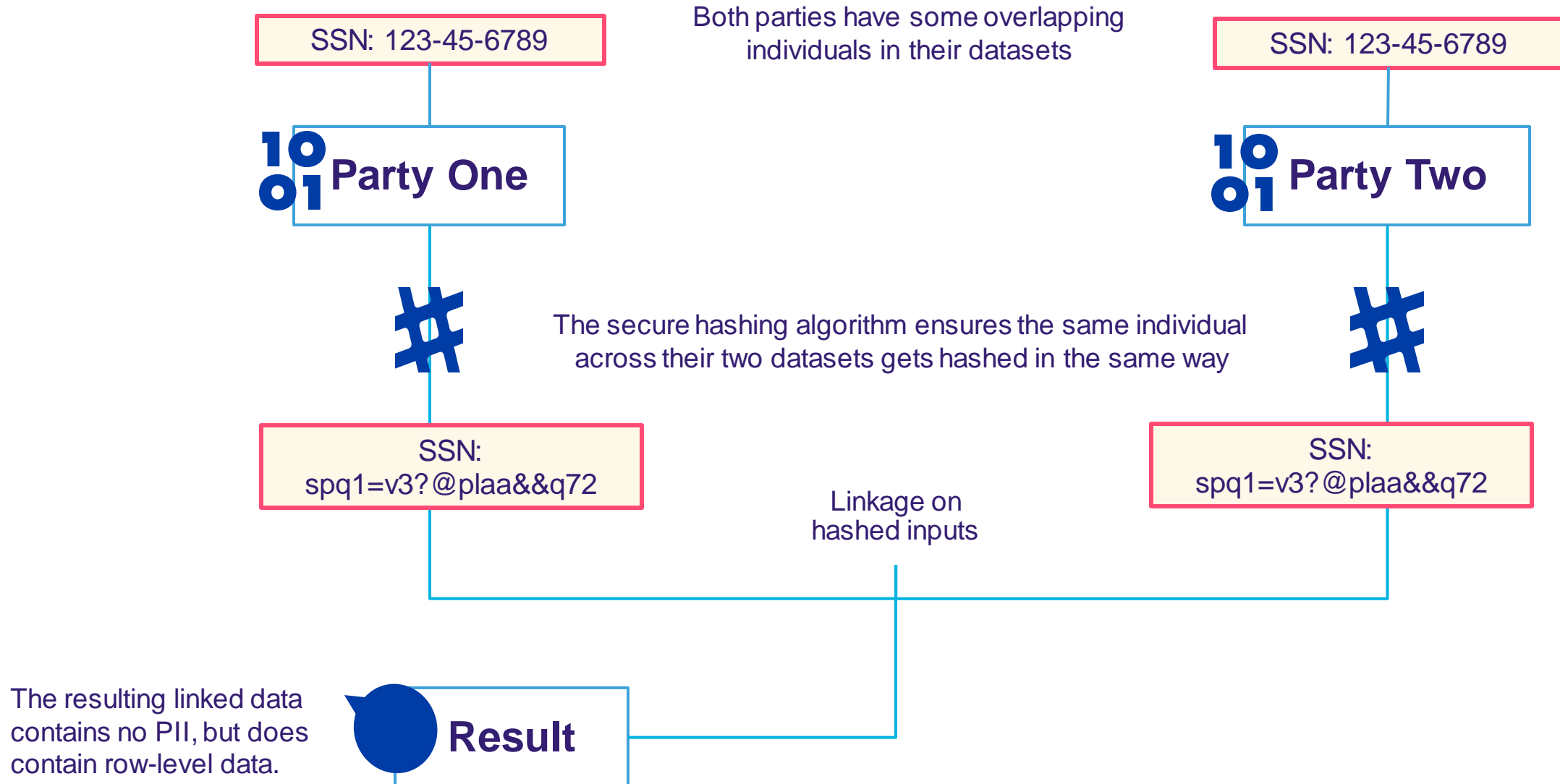
- ☒ generally, not reversible (hash cannot be decrypted, one way)
- ☒ NIST-approved, developed Secure Hashing Algorithms (SHAs)

- ☐ protects identifying or linkage fields though other fields may also be disclosive
- ☐ does not address output privacy

## Examples:

- [Coleridge Initiative's ADRF](#)
- [Georgia Policy Labs](#)
- [CAPriCORN](#) in healthcare
- [N3C system](#) at NCATS
- [ASPE](#) at HHS

# #Secure Hashing





# FAQ: Secure Hashing

Who has access to the secure hashing algorithm?

- Only data owners, or those doing the hashing.

What do you mean it's not reversible?

- Hashing is one-way—there is no 'decryption key' like in encryption. A bad actor would have to find the exact hashing algorithm used and then apply it to the same exact data to figure out which hashes are associated with which individuals.

What is the output from secure hashing?

- Once the sensitive fields, or PII variables, are hashed, they remain that way in the data. Row-level data can be shared.

Why should we trust it?

- Approved by National Institutes of Standards and Technology (NIST), many leverage use of a salt, or random additional data that further complicates the hash, and prevents cryptographic attacks.



# Differential Privacy (DP)

**Differential privacy (DP):** a method for obscuring identities or attributes in the underlying record-level data by infusing data with noise.

- ☒ reduces re-identification risks for individuals or groups in the data (i.e., students, programs)
- ☒ provides a formal privacy guarantee (can guard against threats known today and those in the future)
- ☒ useful for known queries

- ☐ challenging to implement on low levels of geography or unique population groups without adding a lot of noise
- ☐ tradeoff between privacy and accuracy – as you add more “noise” (protection) you move further from true values
- ☐ no input privacy

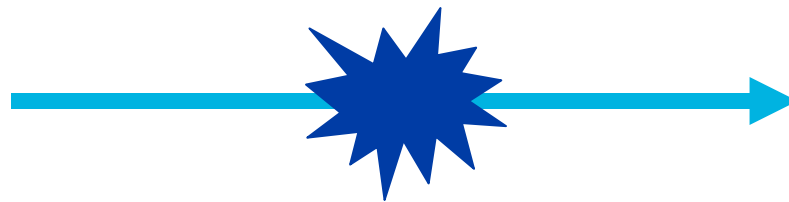
## Examples

- [U.S. Census Bureau's Post-Secondary Employment Outcomes](#)
- [College Scorecard](#)
- [Census 2020](#)
- [Google](#)
- [Apple](#)
- [Facebook](#)



# Differential Privacy (DP)

10  
01



10  
01

A County	323
B County	4,002
D County	5,275
E County	427
F County	112
G County	3,936
H County	230
I County	75

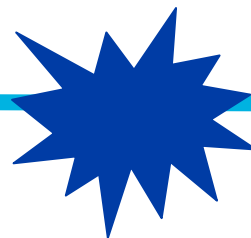
+	27
-	24
-	35
+	16
+	14
-	17
+	8
+	12

A County	350
B County	3,978
D County	5,240
E County	443
F County	126
G County	3,919
H County	238
I County	87



# Differential Privacy

(non-numeric fields)



Student A	White
Student B	Black
Student C	Asian
Student D	White
Student E	AIAN
Student F	White
Student G	Black
Student H	White
Student I	Asian
Student J	White
Student K	White
Student L	White

...

...

change  
don't change  
change  
don't change  
don't change  
don't change  
don't change  
change  
don't change  
don't change  
change  
don't change

coin  
toss

Student A	Asian
Student B	Black
Student C	White
Student D	White
Student E	AIAN
Student F	White
Student G	Black
Student H	Asian
Student I	Asian
Student J	White
Student K	White
Student L	White

...

...

# FAQ: Differential Privacy

Who gets to decide how much privacy vs. utility to preserve?

- The data owners. They can work with DP experts to tweak the mathematical formula to their data needs.

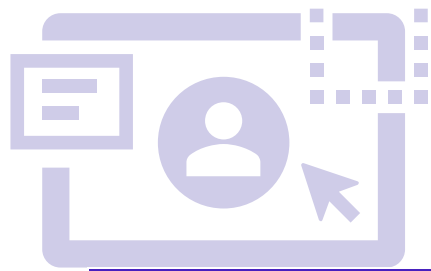
Can DP only apply to data after analysis is done?

- No, you could apply DP to data BEFORE analysis is done (local) or AFTER aggregate statistics have been produced (central/global).

Why do I need this? Isn't swapping or aggregation enough?

- Computing power has grown in recent decades so that individual bad actors can quickly and easily re-trace the identities of people in a dataset. No other disclosure limitation protocol provides a mathematical guarantee of privacy in the way that DP does.

# What it takes to implement



Technical



Legal



Cultural



Institutional

# Questions?

# Share your thoughts

Take a quick  
[Workshop Survey](#)

For more trainings, visit:  
<https://disc.wested.org/>







# Thank you.

---

**Presenter 1**

Title  
email

**Presenter 1**

Title  
email

A Project of  
WestEd 



Copyright ©2024 Data Integration Support Center at WestEd and Actionable Intelligence for Social Policy at University of Pennsylvania.