

# SYNTHETIC DATA GENERATION

**FEB. 9, 2022 | 11AM EDT**

# 101

Presented by



Dr. Khaled El Emam,  
SVP and General Manager,  
Replica Analytics

# Synthetic Data Generation 101

*Khaled El Emam*

*February 9, 2022*



AN AETION COMPANY

# Agenda

## Introduction to Synthesis

1

General description of what synthetic data is and general use cases

## Privacy and Utility

2

An examination of privacy risks and the utility of synthetic data

## FAQs

3

Some commonly asked questions about synthetic datasets

## Additional Use Cases

4

Beyond data sharing, SDG has additional applications



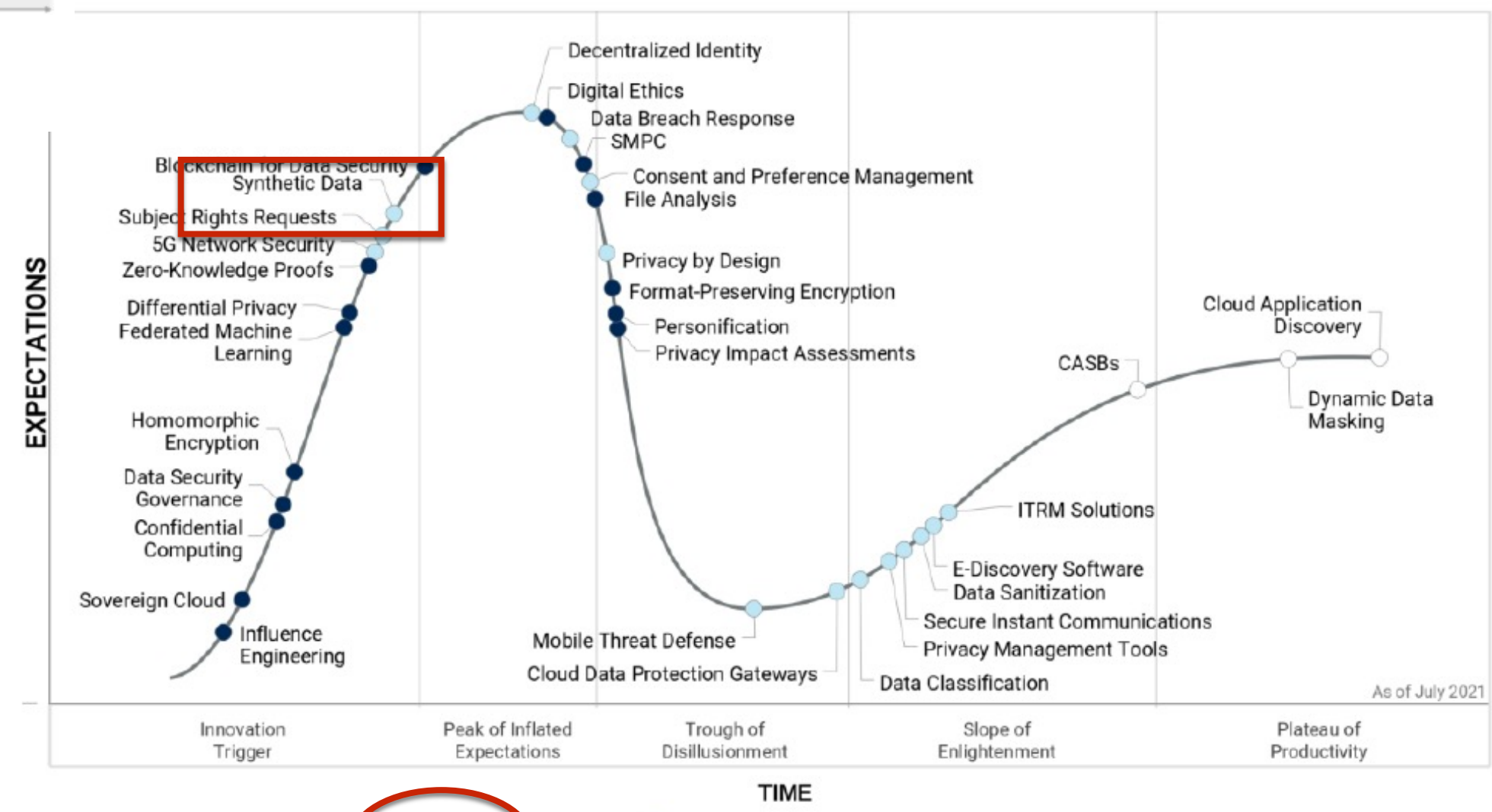
# The adoption of synthetic data has been accelerating quite rapidly



2020

2021

Plateau will be reached:  
 ○ less than 2 years   ● 2 to 5 years   ● 5 to 10 years   ▲ more than 10 years   ✗ obsolete before plateau



Plateau will be reached: ○ < 2 yrs.   ● 2-5 yrs.   ● 5-10 yrs.   ▲ >10 yrs.   ✗ Obsolete before plateau

Gartner  
Hype Cycle for Privacy, 2021

# Gartner predicts synthetic data will have a non-trivial impact on privacy violations and sanctions

## Top 10 Strategic Predictions for 2022 and Beyond

Data	Tracking	Behavior	Supervision	Talent
<b>70%</b> reduction in privacy sanctions	<b>40%</b> intentionally devalue personal data	<b>25%</b> neuromine at scale	<b>30%</b> teams without a boss	<b>30%</b> increase in talent across Africa
2025	2024	2027	2024	2026
Composability	Cyber Attack	Customers	Crypto	Digital
<b>80%</b> report better business performance	<b>G20</b> cyber attack breeds kinetic response	<b>75%</b> companies "break up" with customers	<b>NFTs</b> drive high value companies	<b>1 Billion</b> poorest people get internet
2024	2024	2025	2026	2027

[gartner.com](https://gartner.com)

Source: Gartner  
© 2021 Gartner, Inc. and/or its affiliates. All rights reserved. CTMKT\_1544427

**Gartner**



AN AETION COMPANY

# The Erosion of Trust ?

The New York Times

## *Your Data Were 'Anonymized'? These Scientists Can Still Identify You*

Computer scientists have developed an algorithm that can pick out almost any American in databases supposedly stripped of personal information.

Opinion | [THE PRIVACY PROJECT](#)

## Twelve Million Phones, One Dataset, Zero Privacy

By Stuart A. Thompson and Charlie Warzel  
DEC. 19, 2019

ACM TECHNEWS

## 'Anonymized' Data Can Never Be Totally Anonymous, says Study

By The Guardian

## Online Profiling and Invasion of Privacy: The Myth of Anonymization

02/20/2013 12:23 pm ET | Updated Apr 22, 2013

theguardian

## 'Anonymised' data can never be totally anonymous, says study

Findings say it is impossible for researchers to fully protect real identities in datasets

## You're very easy to track down, even when your data has been anonymized

A new study shows you can be easily re-identified from almost any database, even when your personal details have been stripped out.

by Charlotte Jee

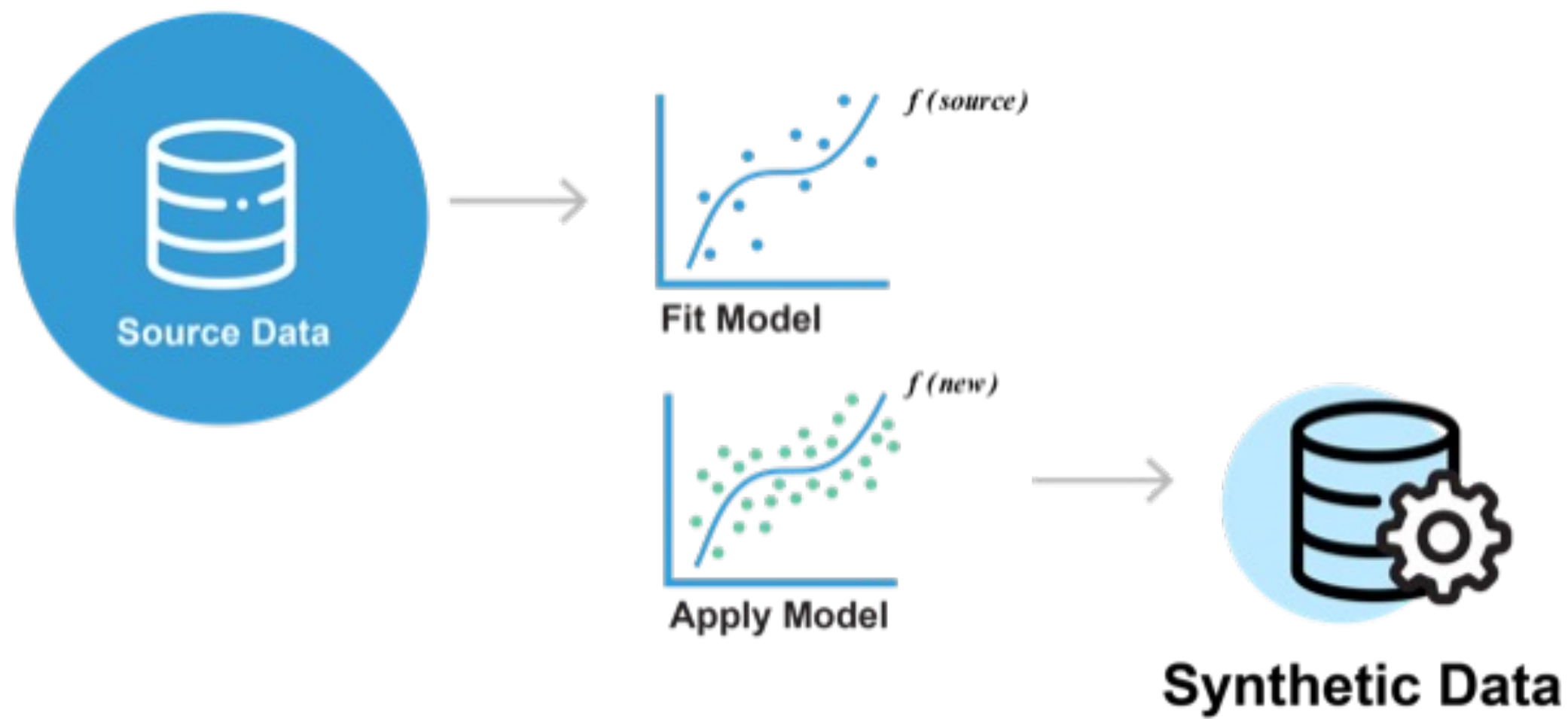
Jul 23, 2019

HUFFPOST



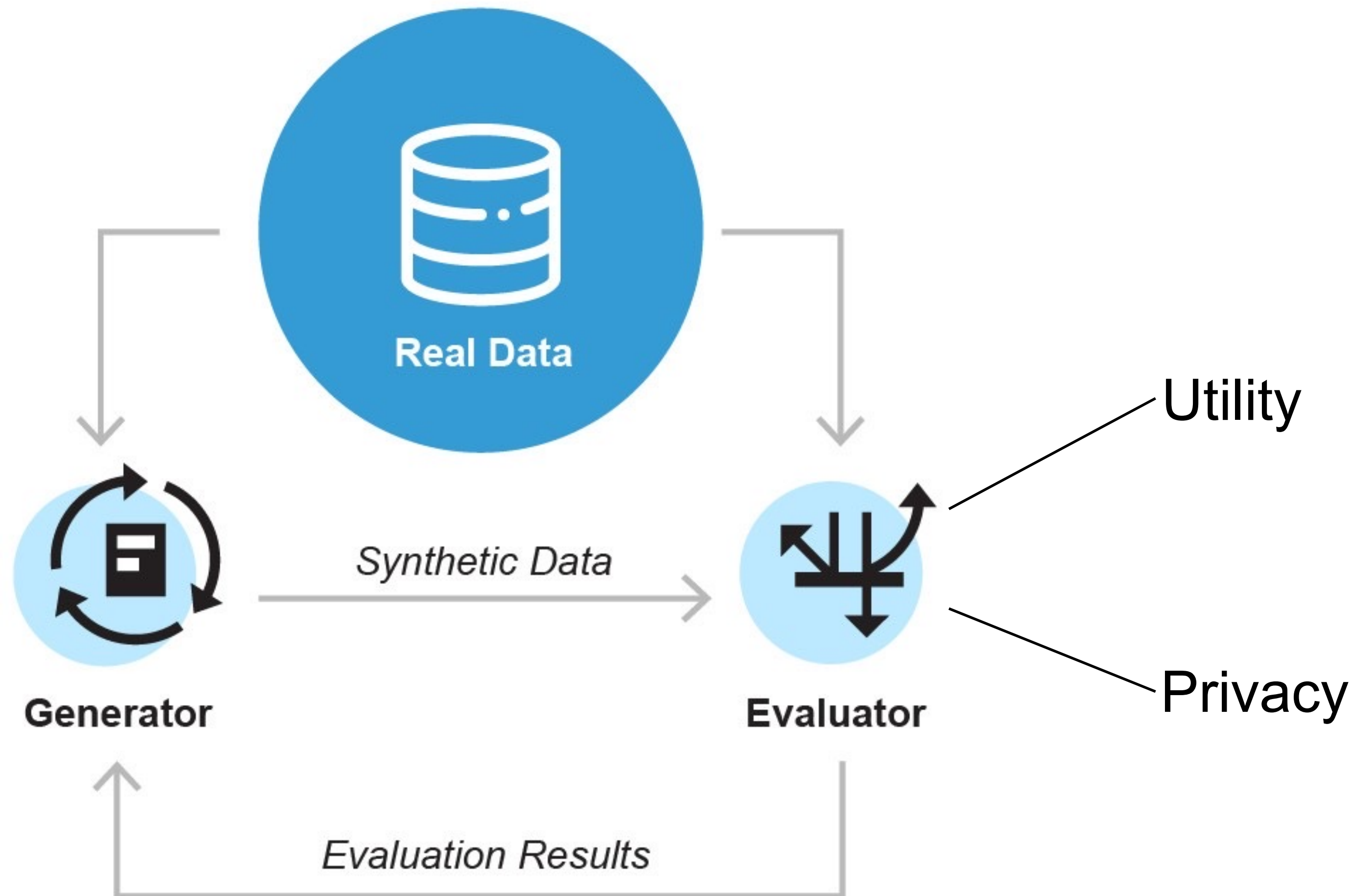
AN AETION COMPANY

# The Synthesis Process



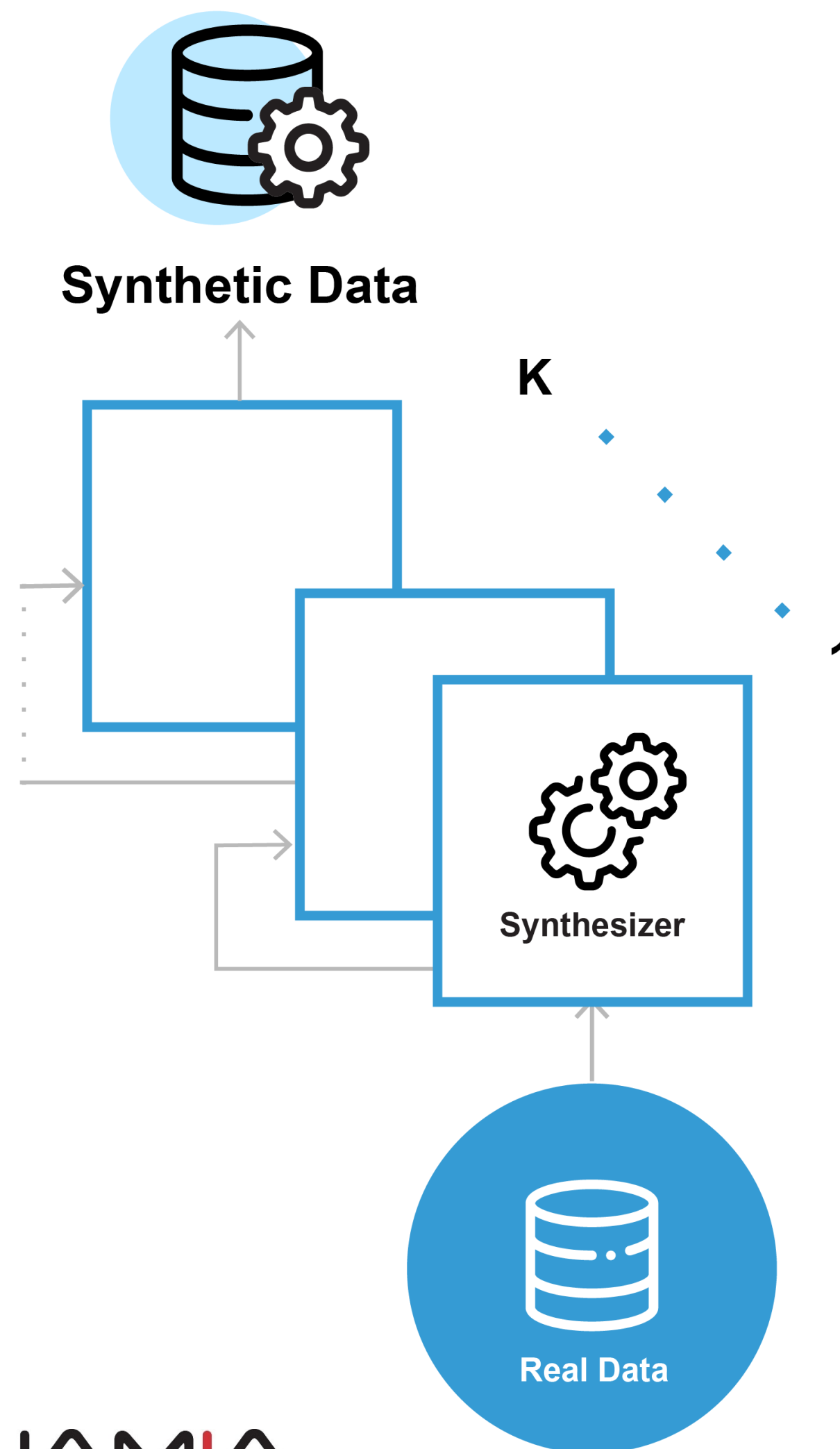
COU1A	AGECAT	AGELE70	WHITE	MALE	BMI
United States	2	1	1	1	33.75155
United States	2	1	1	0	39.24707
United States	1	1	1	0	26.5625
United States	4	1	1	1	40.58273
United States	5	0	0	1	24.42046
United States	5	0	1	0	19.07124
United States	3	1	1	1	26.04938
United States	4	1	1	1	25.46939

# Training a generative model uses a utility – privacy loss function



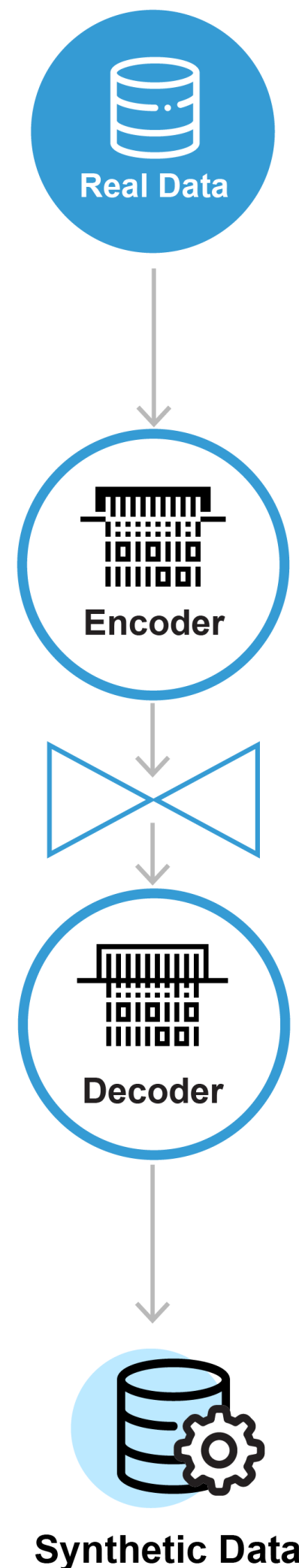


# Sequential synthesis utilizes multiple machine learning methods in a sequence

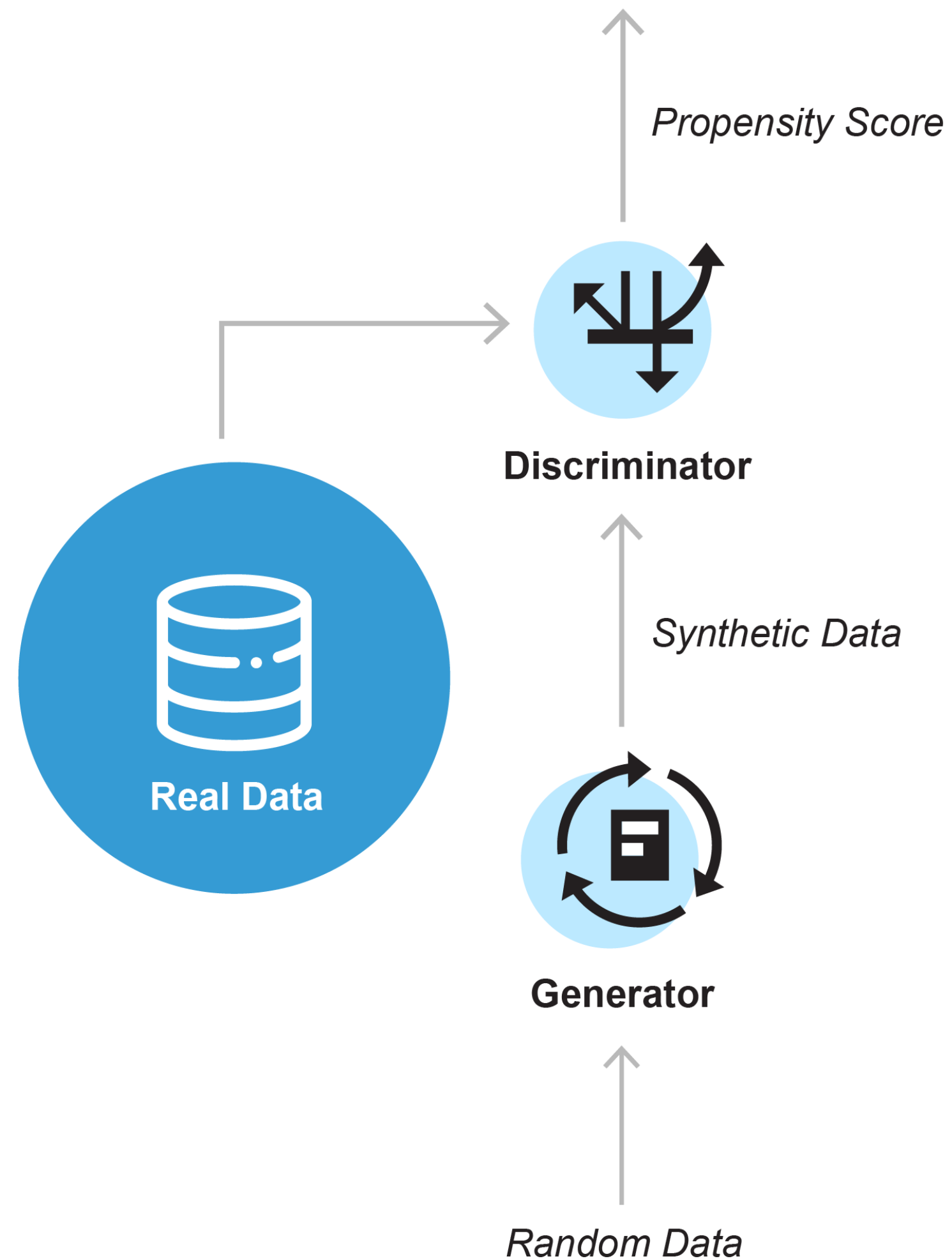


# Variational Auto Encoder (VAE)

compresses the input into a latent space



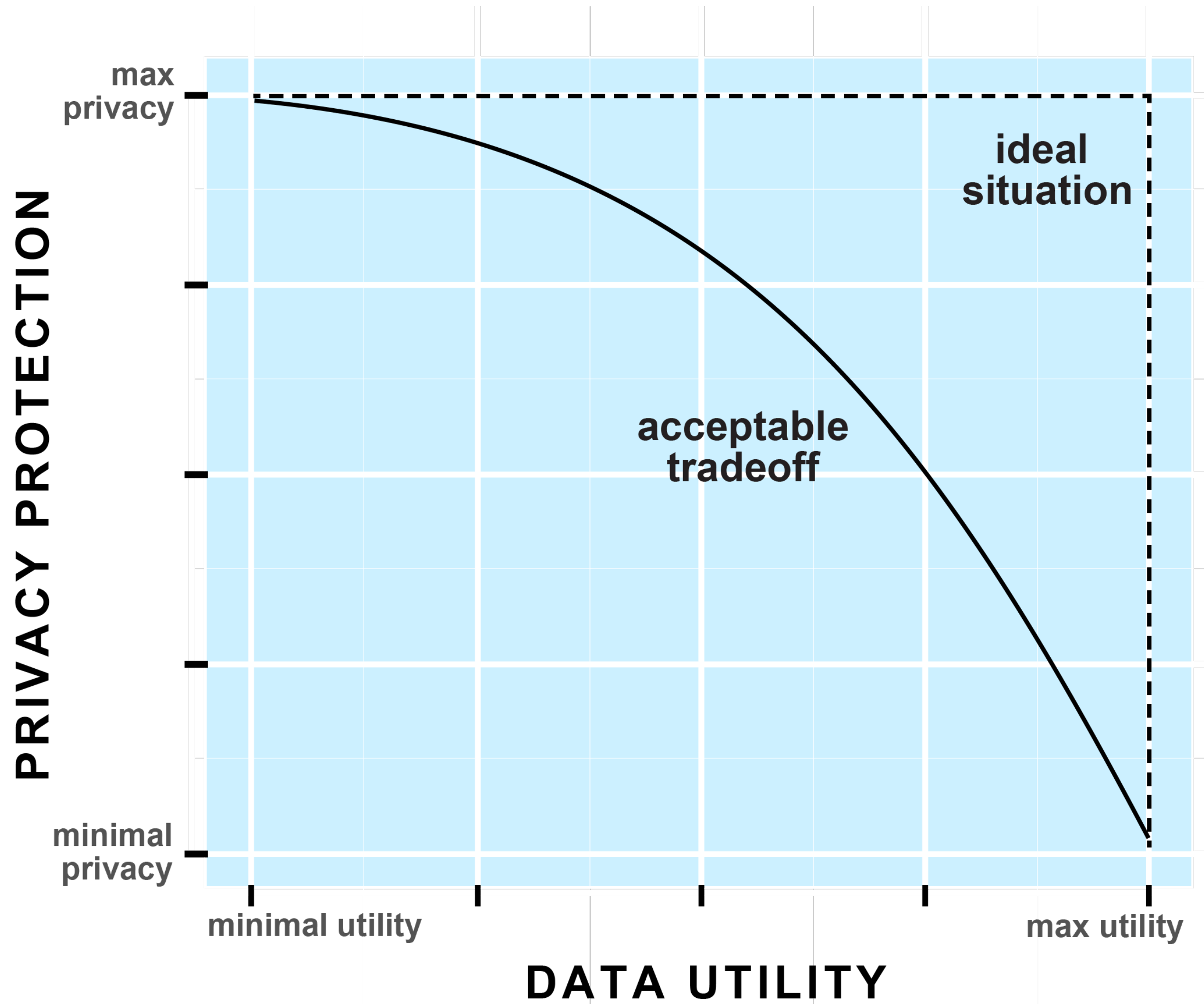
# Generative Adversarial Network (GAN)



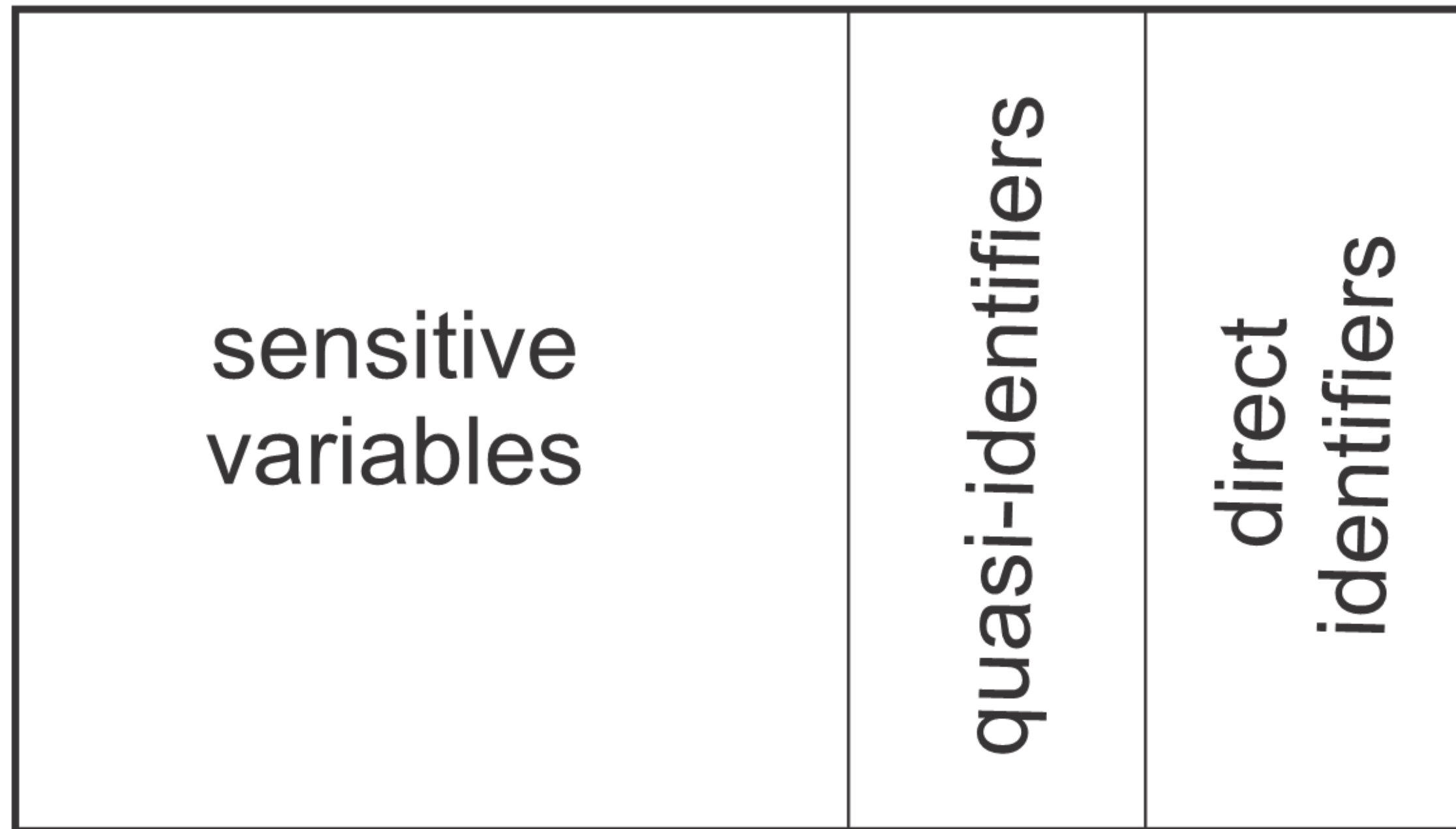
# There are seven common use cases for synthetic data

1. Machine learning  
*(model evaluation, data augmentation, sharing ML models)*
2. Software testing
3. Education, training, and hackathons
4. Data retention
5. Vendor assessment
6. Internal secondary use  
*(exploratory and detailed analytics)*
7. External data sharing

# Privacy-Utility Trade-off

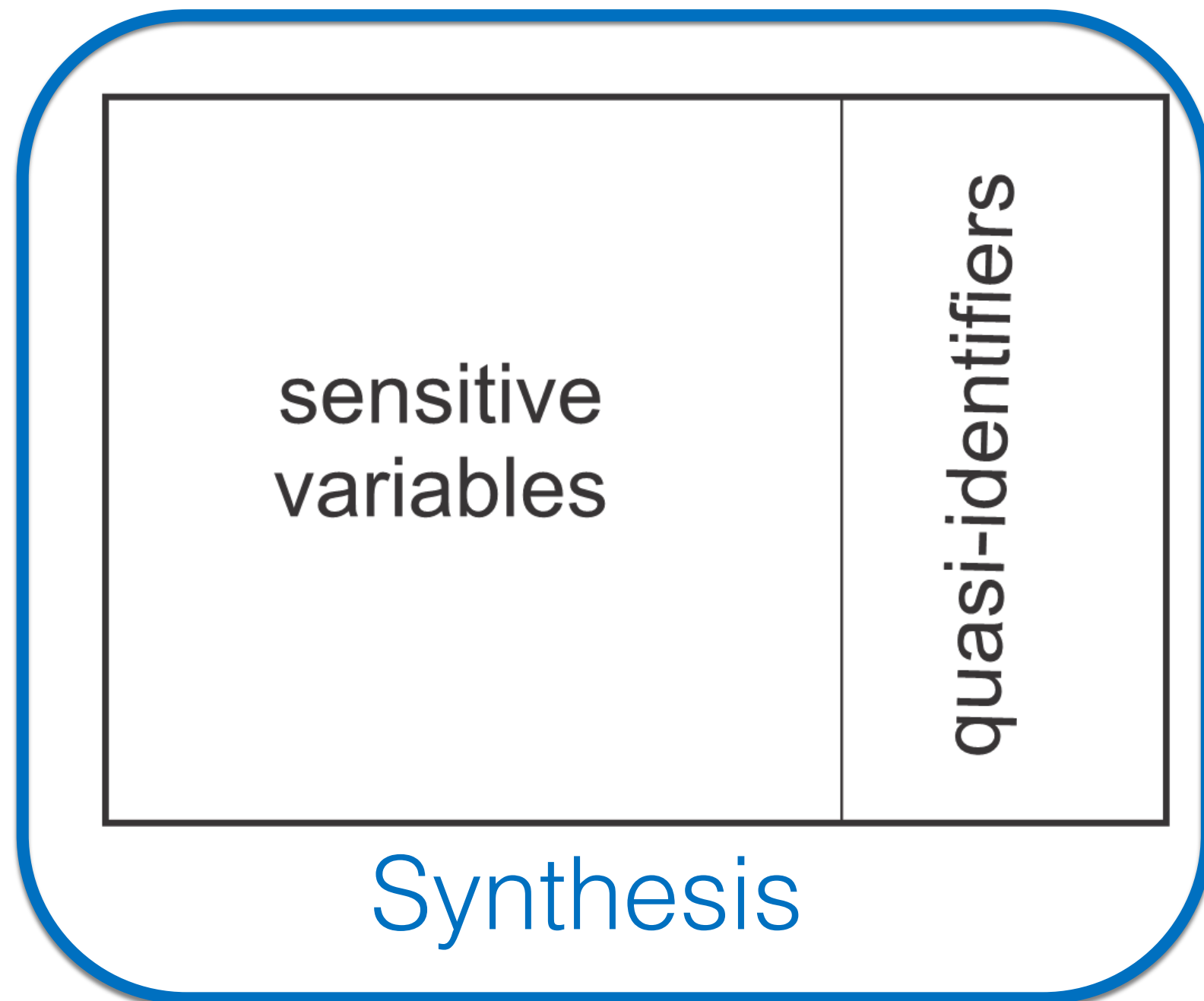


# Variables in a dataset can be classified into one of three types

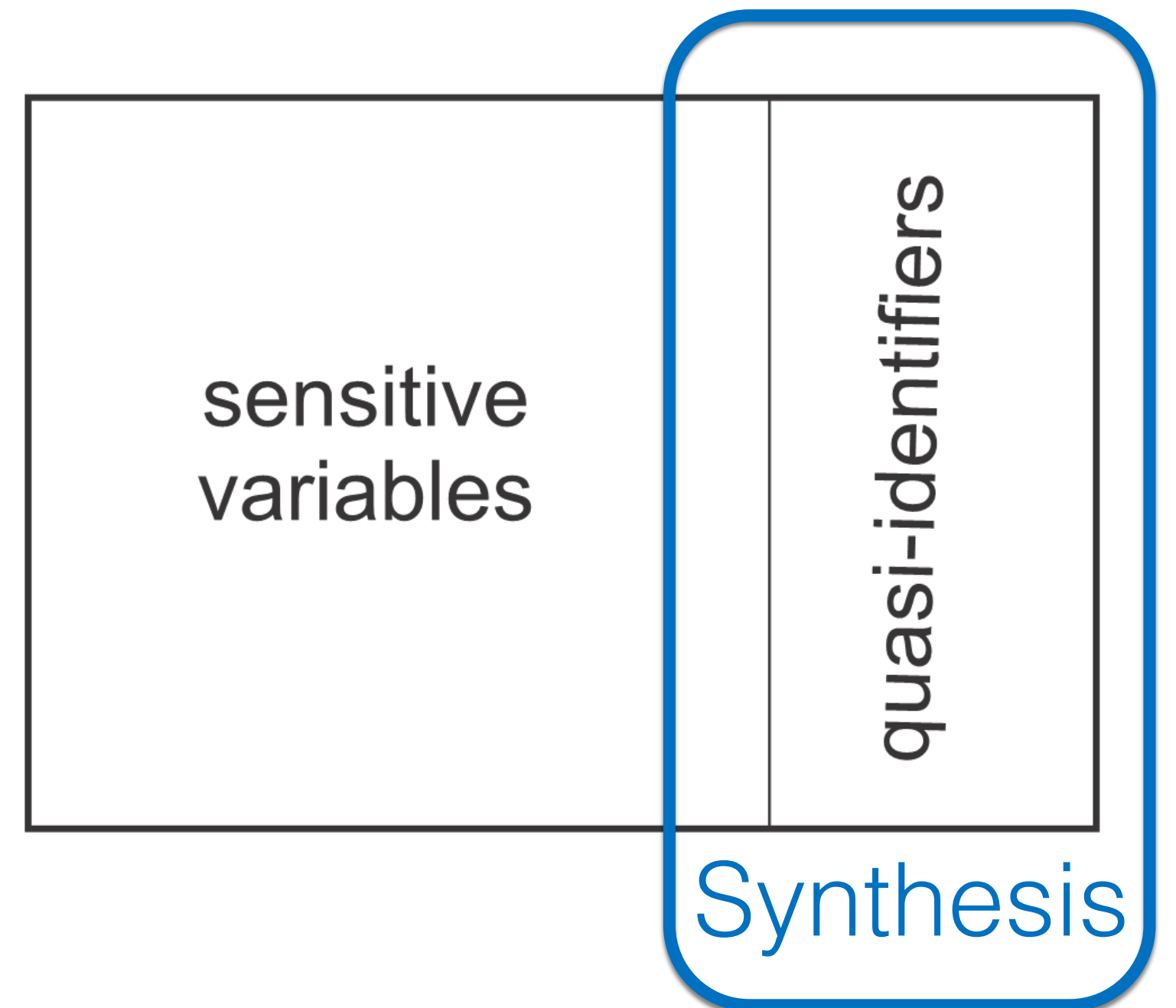


# Two Synthesis Strategies

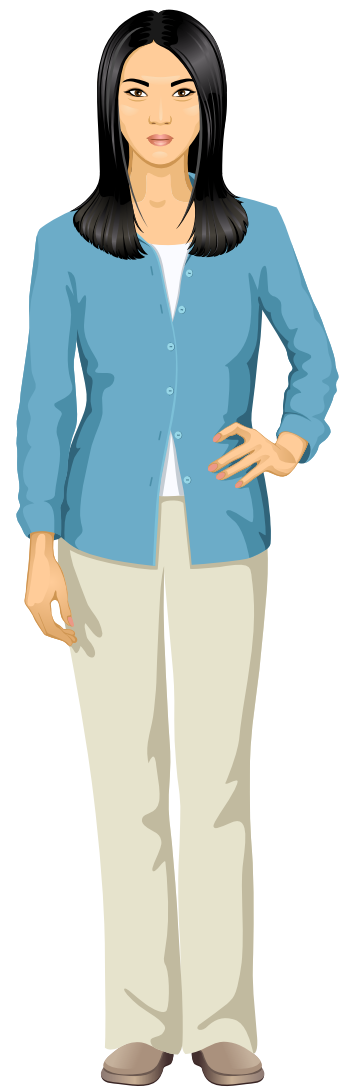
Full Synthesis  
Synthesize all  
variables



Partial Synthesis  
Synthesize  
quasi-identifiers



# Attribution disclosure: find a similar record in the synthetic data and learn something new



Quasi-identifiers



New Information



Sex	Year of Birth	NDC
Male	1975	009-0031
Male	1988	0023-3670
Male	1972	0074-5182
Female	1993	0078-0379
<b>Female</b>	<b>1989</b>	<b>65862-403</b>
Male	1991	55714-4446
Male	1992	55714-4402
Female	1987	55566-2110
Male	1971	55289-324
Female	1996	54868-6348
Male	1980	53808-0540

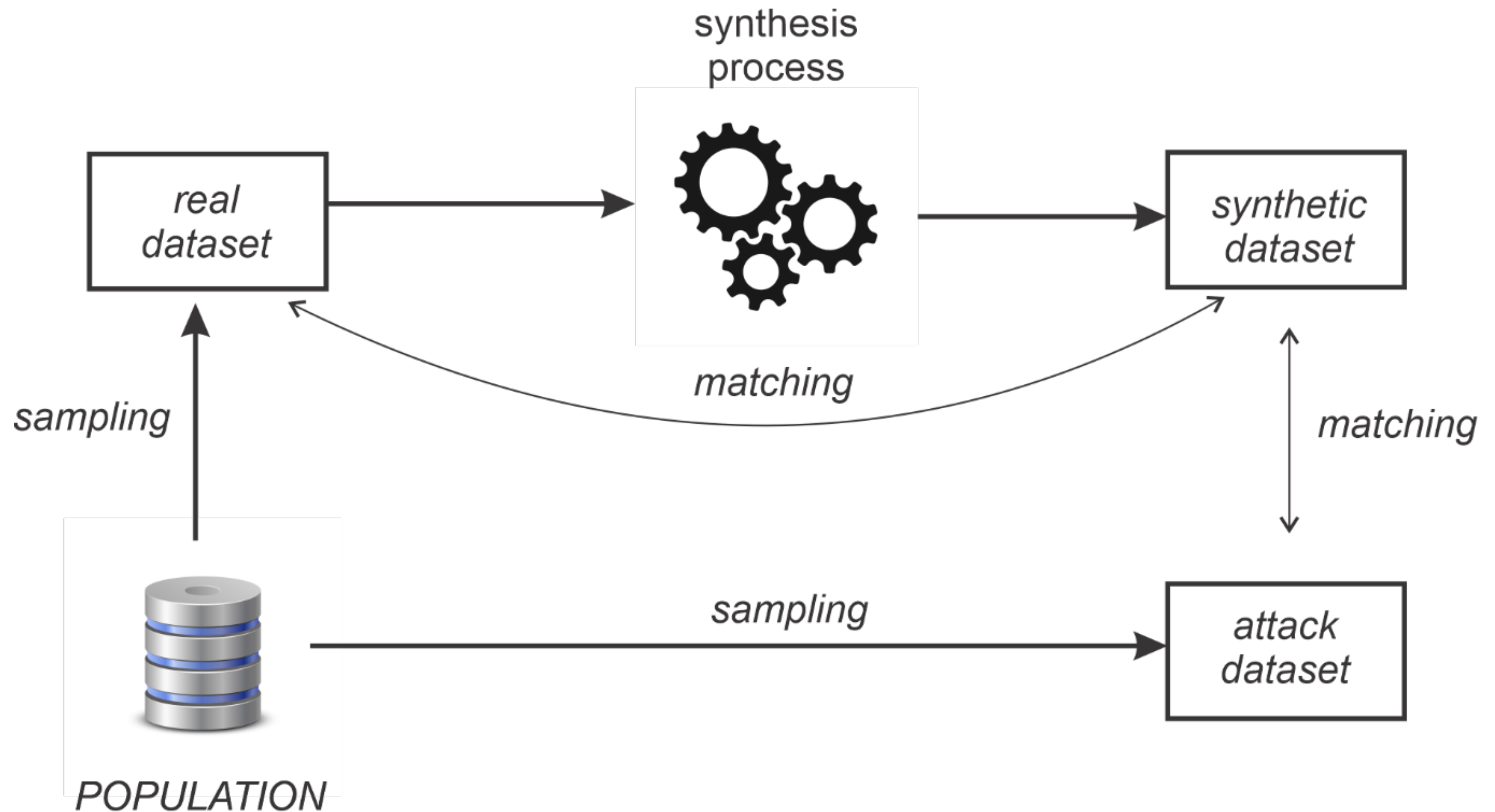


# Evaluations of attribution risks show that it is low in multiple studies across multiple datasets

Dataset	Fully Synthetic Data	Original Data
Washington Hospital Data (Discharge)	0.0197	0.098
Canadian COVID-19 Data (Public Health)	0.0086	0.034

A commonly used risk threshold = 0.09

# Membership disclosure



# One way to classify utility metrics is as broad and narrow

## broad metrics → narrow metrics

These are generic metrics that are easy to calculate when the generative model is built and synthetic data are synthesized. They are only useful if they are predictive of workload-specific metrics.

There are multiple use cases for these metrics:

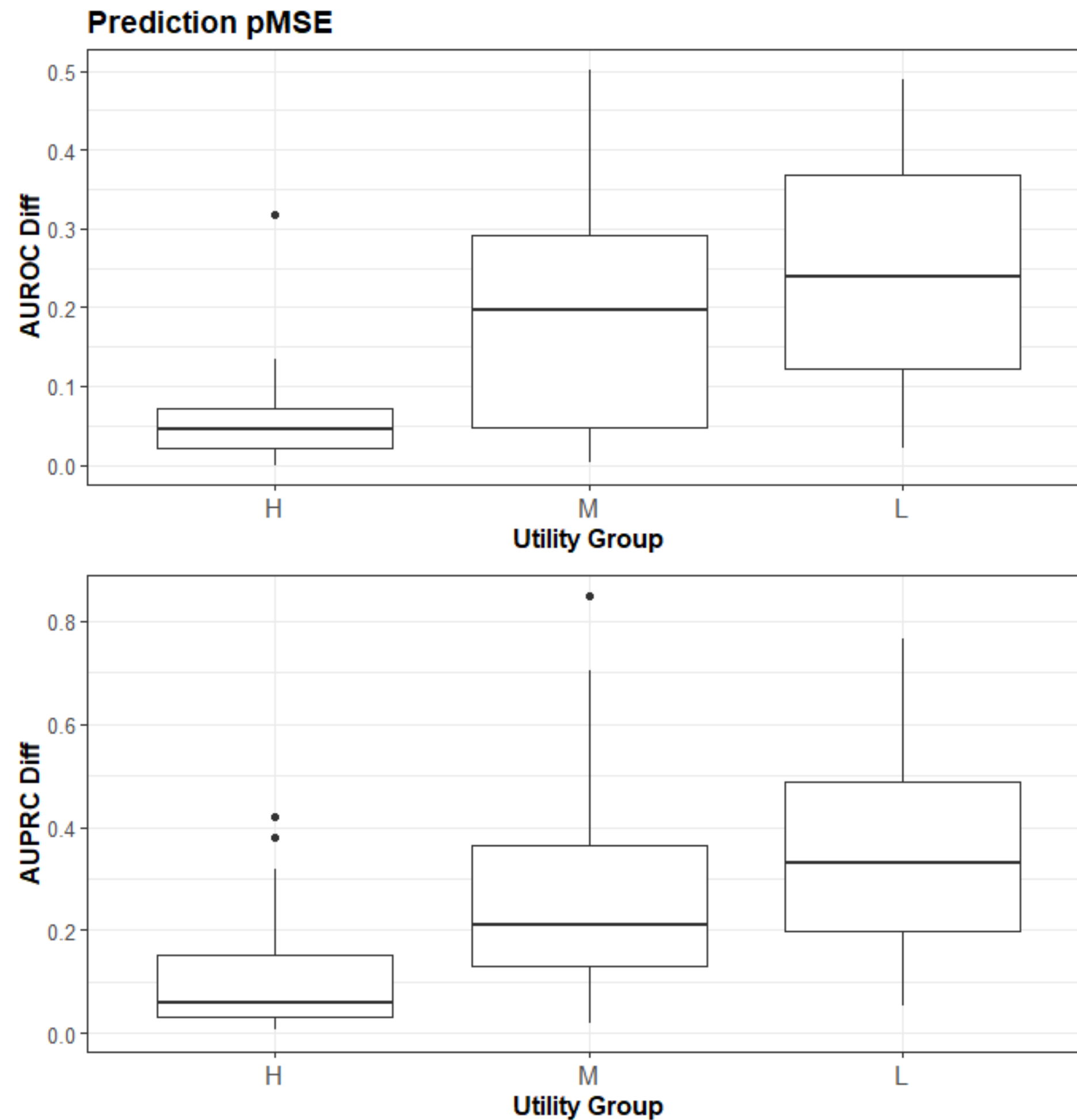
- Dataset vs model
- Use case: rank / hyperparameter tuning / communicate

These are workload-specific and are what is of most interest to the data users. However, all the possible workloads will not be known in advance and therefore we have to consider representative workloads when developing and evaluating utility metrics.

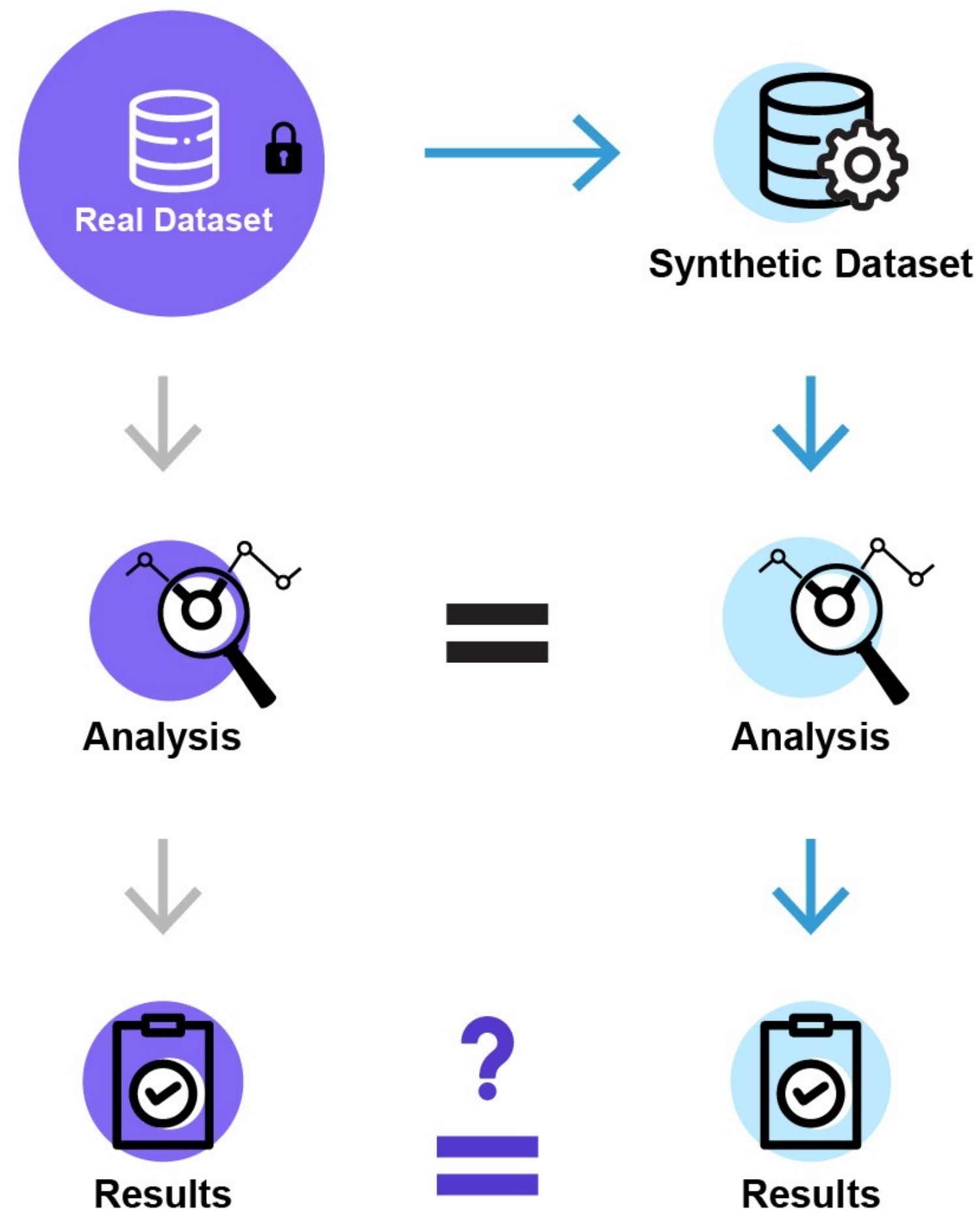
There are different types:

- Information loss metrics
- Inferential validity metrics

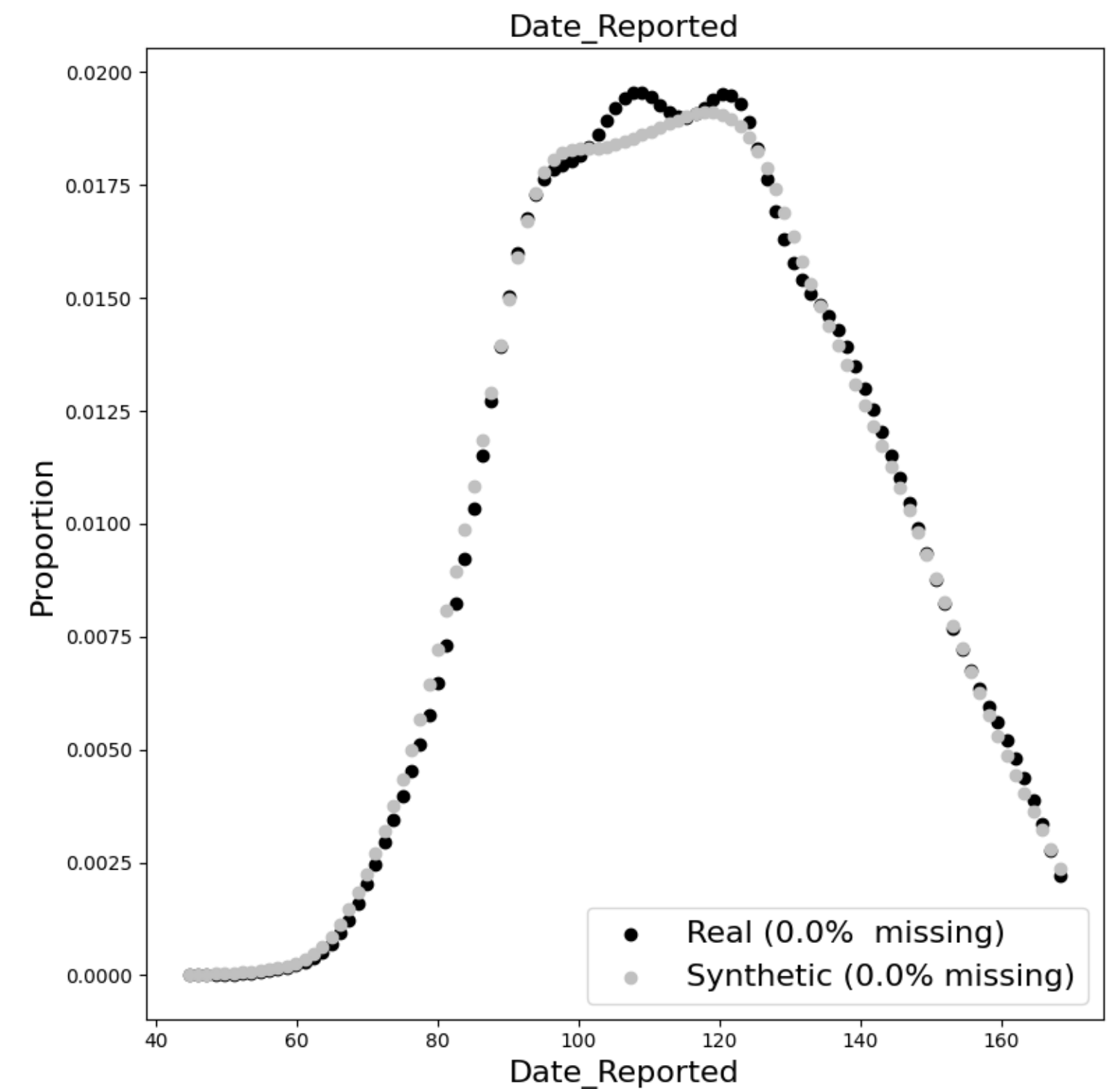
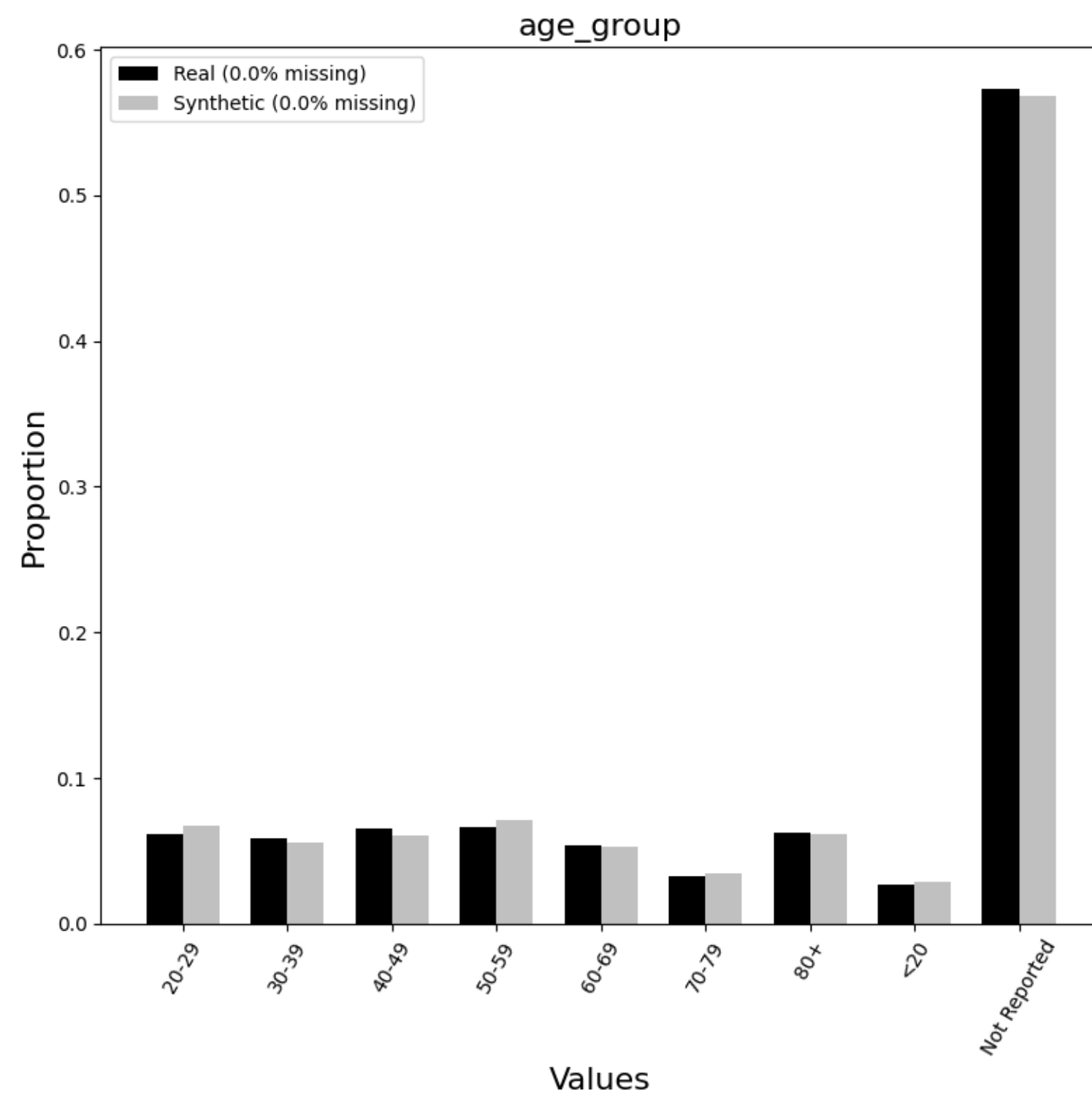
# Broad utility metrics can rank SDG methods by their workload performance



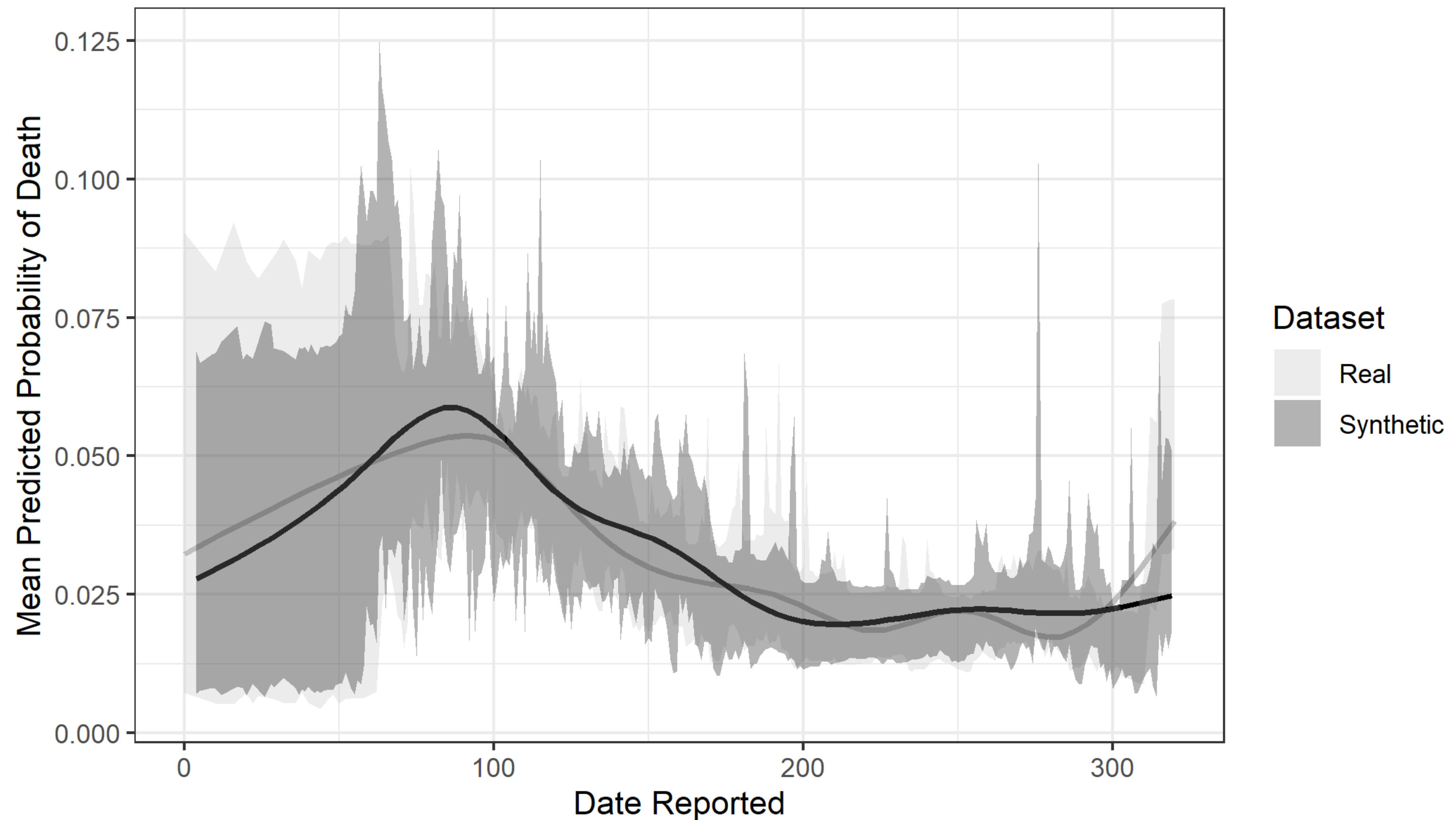
# To evaluate utility one can compare the analysis results from real and synthetic data



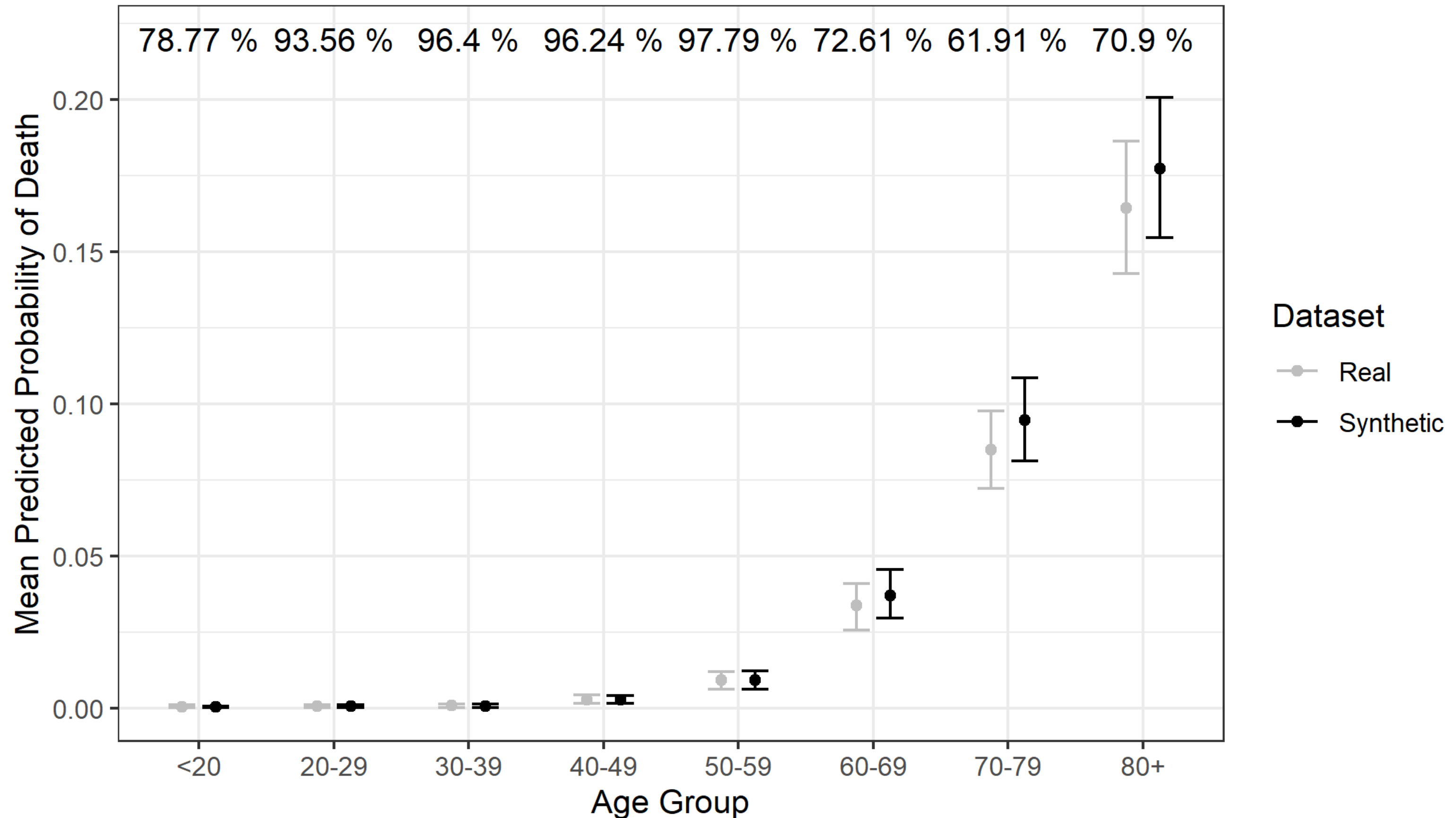
# The univariate distributions of real and synthetic datasets look similar



# Mortality over time for the Ontario COVID-19 case dataset



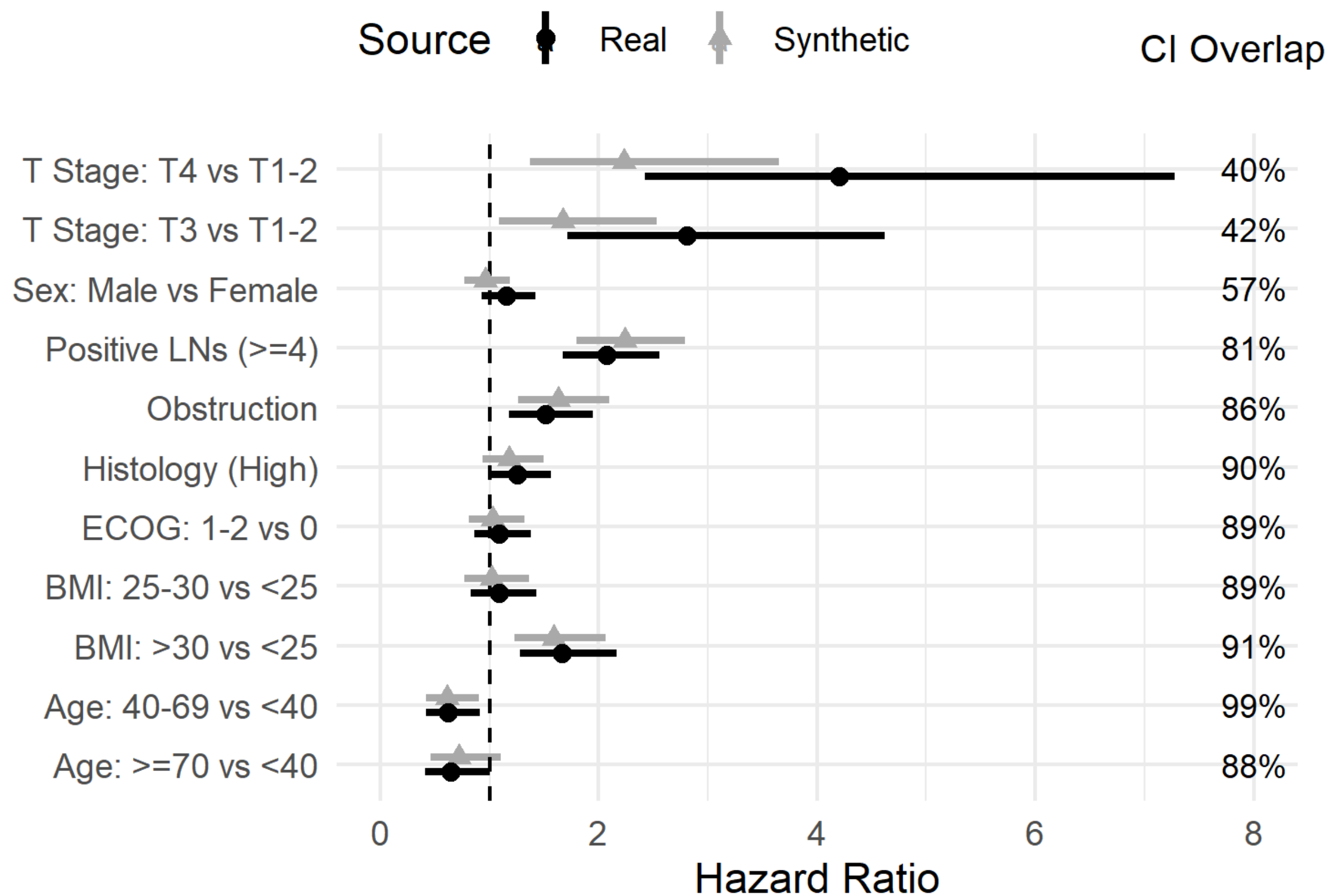
# Mortality by age for the Ontario COVID-19 case dataset



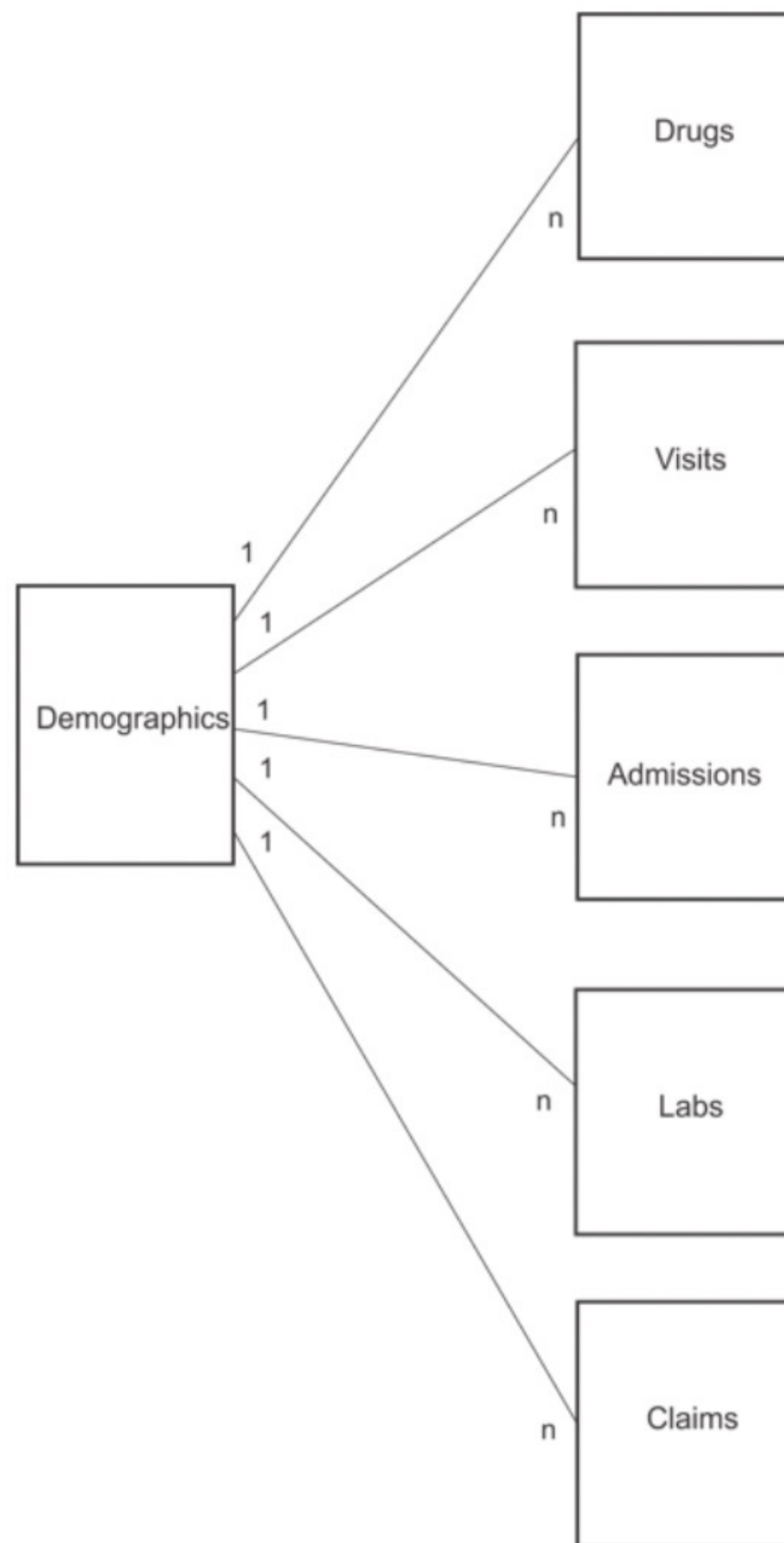


# Comparing real and synthetic data: Adjusted model of impact of bowel obstruction on DFS

## Hazard Ratios: Analysis for Disease-Free Survival



# Longitudinal Data Model



Demographics
Age
Sex
Time to last day of follow-up available
Comorbidity score (elixhauser)

Drugs
Dispensed amount quantity
Relative dispensed time in days
Dispensed day supply quantity
Morphine use (binary)
Oxycodone use (binary)
Antidepressant use (binary)

Visits (ED)
Relative admission time in days
Problem code 1
Problem code 2
Resource intensity weights

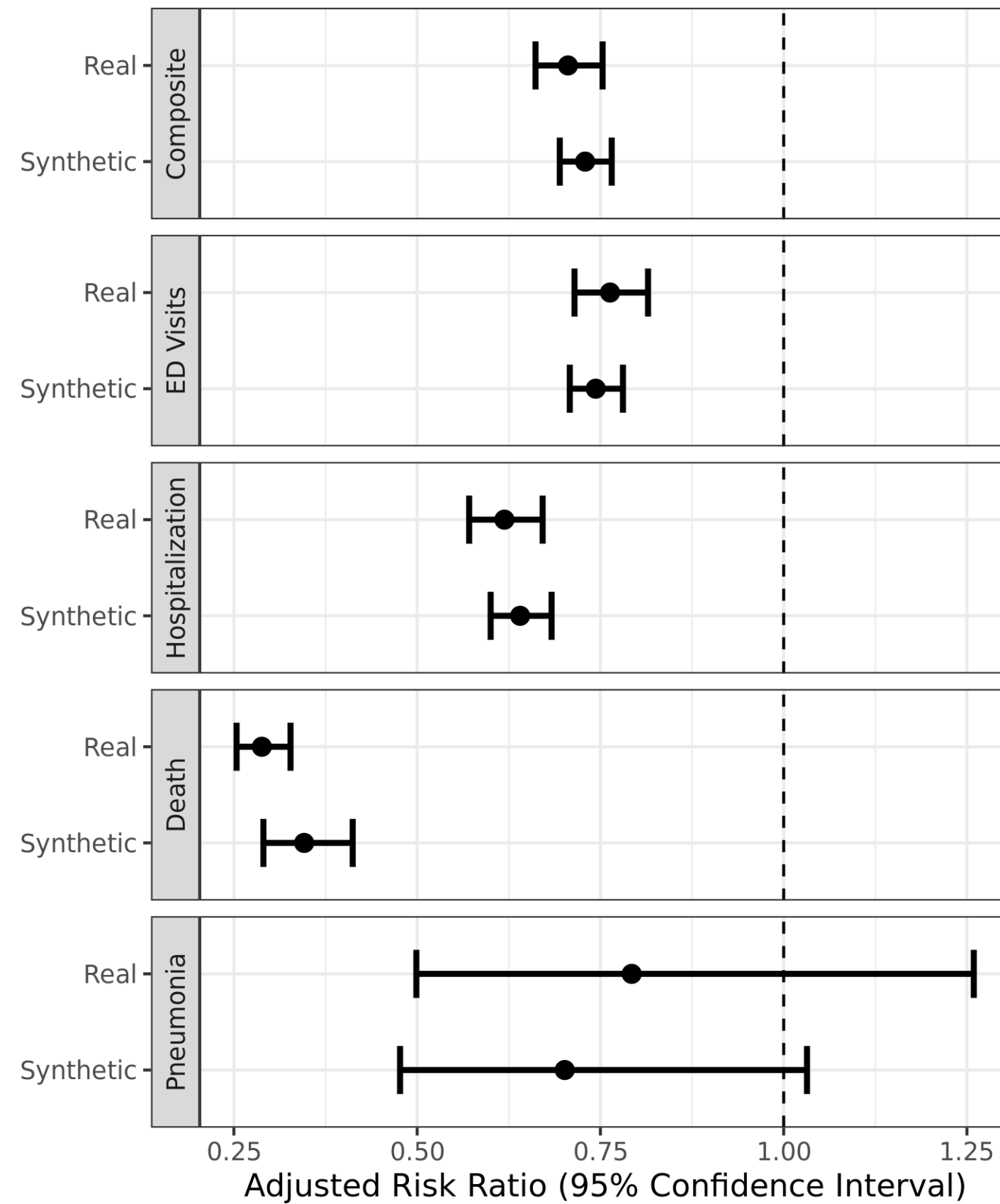
Admissions (Hospital)
Relative time admitted in days
LOS
Diagnosis code 1
Diagnosis code 2
Resource intensity weight

Lab
Test name
Test result (integer)
Relative time in days lab taken

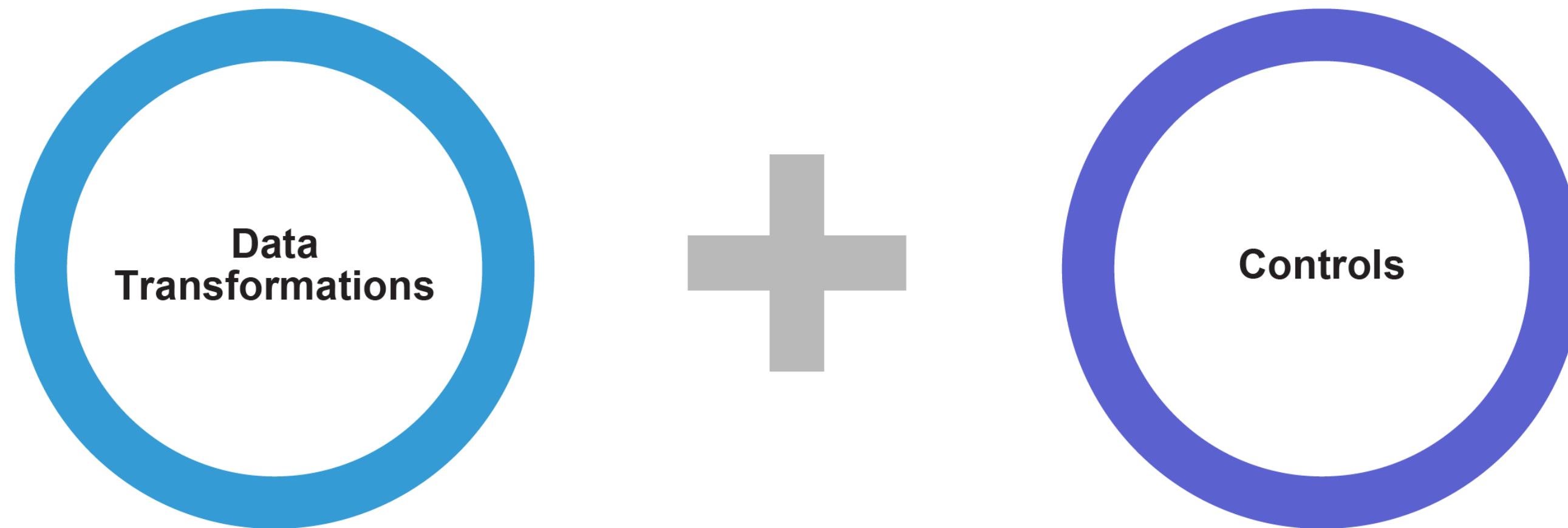
Claims
Primary diagnosis code
Provide specialty
Relative service event start date

# Adjusted Cox Regression

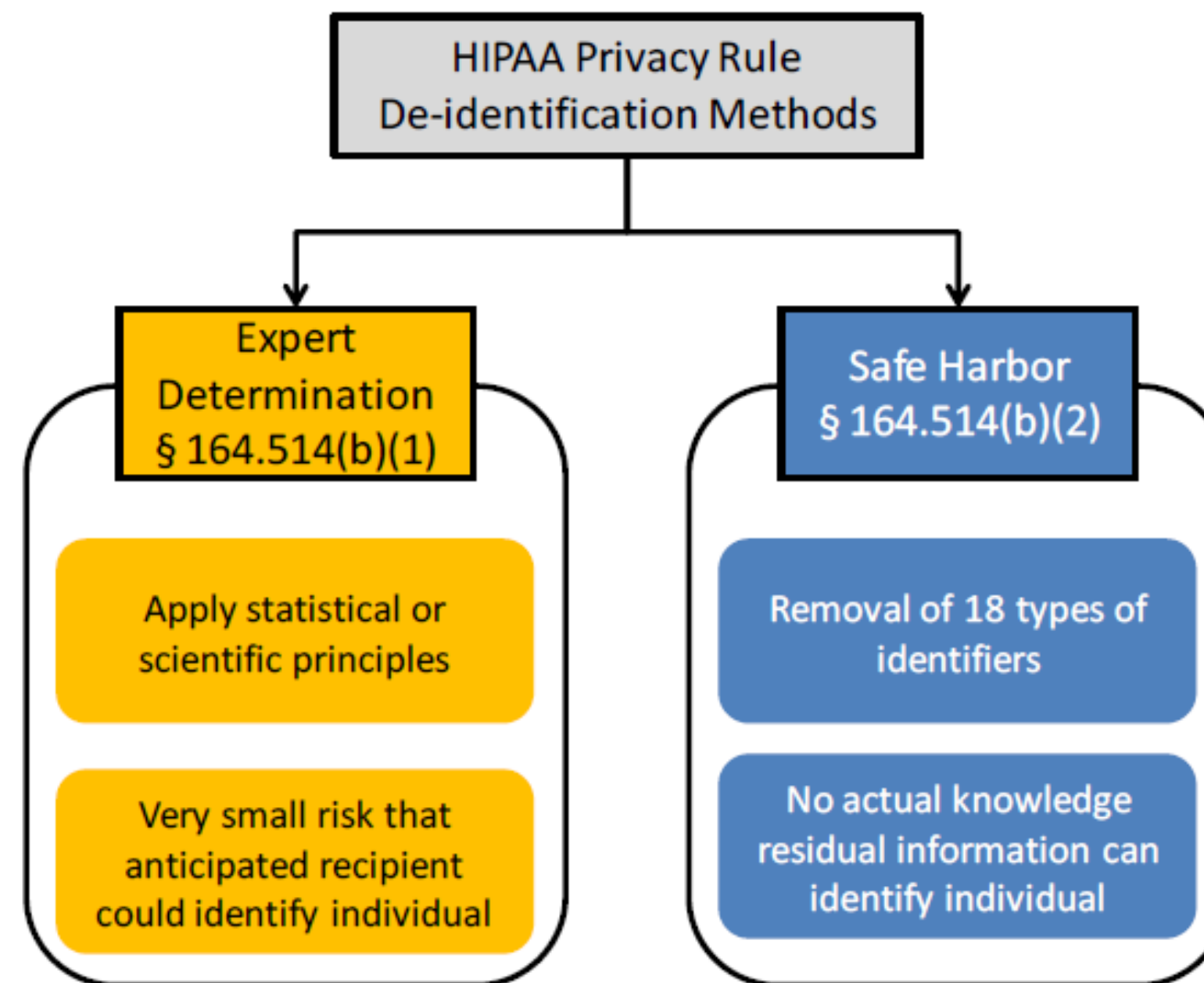
Note: Adjusted estimates include the following co-variates: age, sex, antidepressant use, Elixhauser score, ALT, eGFR, HCT; Opioid 1 served as the reference group



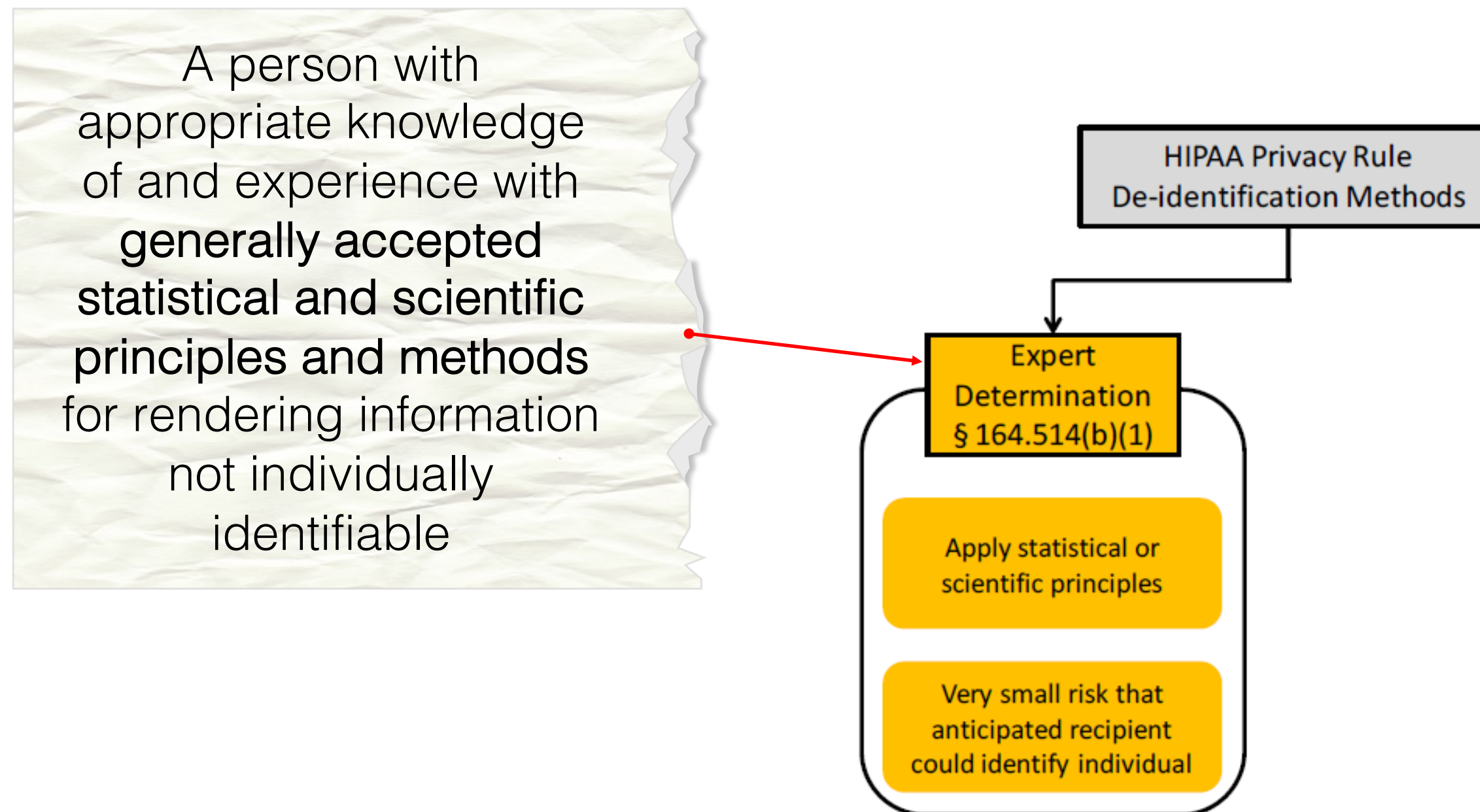
# Are controls needed to manage the privacy risks in synthetic data ?



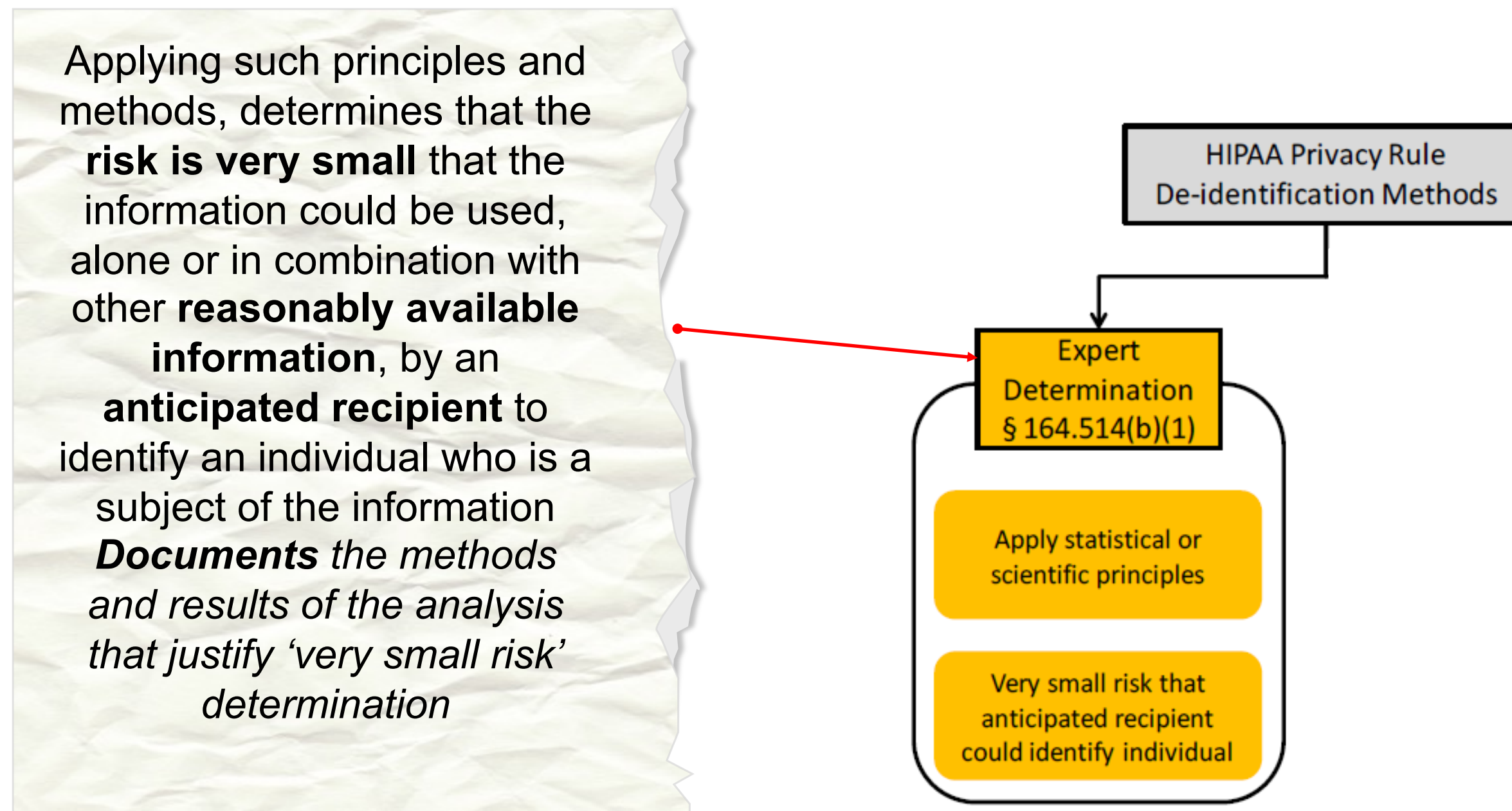
# HIPAA de-identification methods are still being used – they are generally considered the gold standard



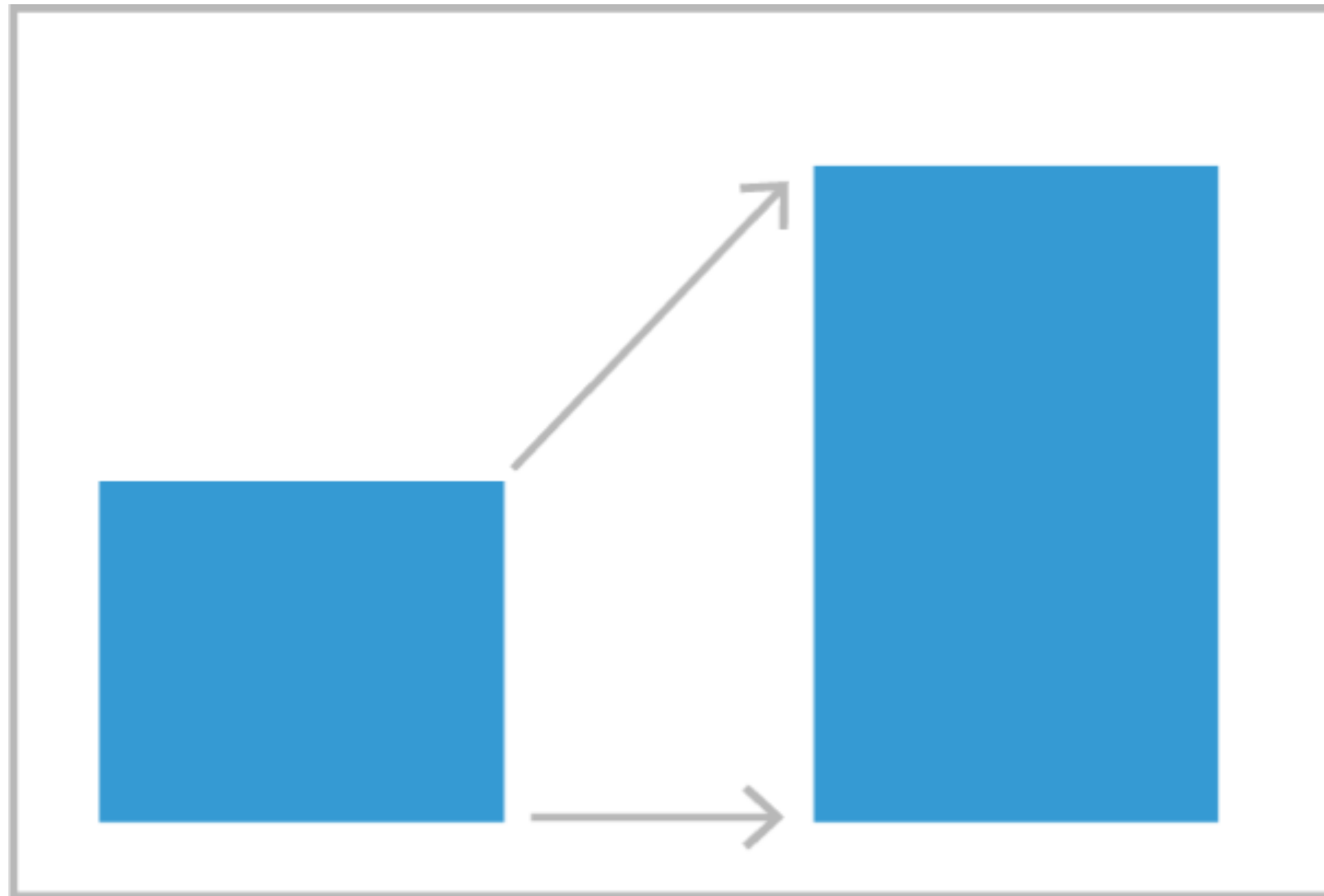
# The Expert Determination Method



# The Expert Determination method is more consistent with modern disclosure control practices

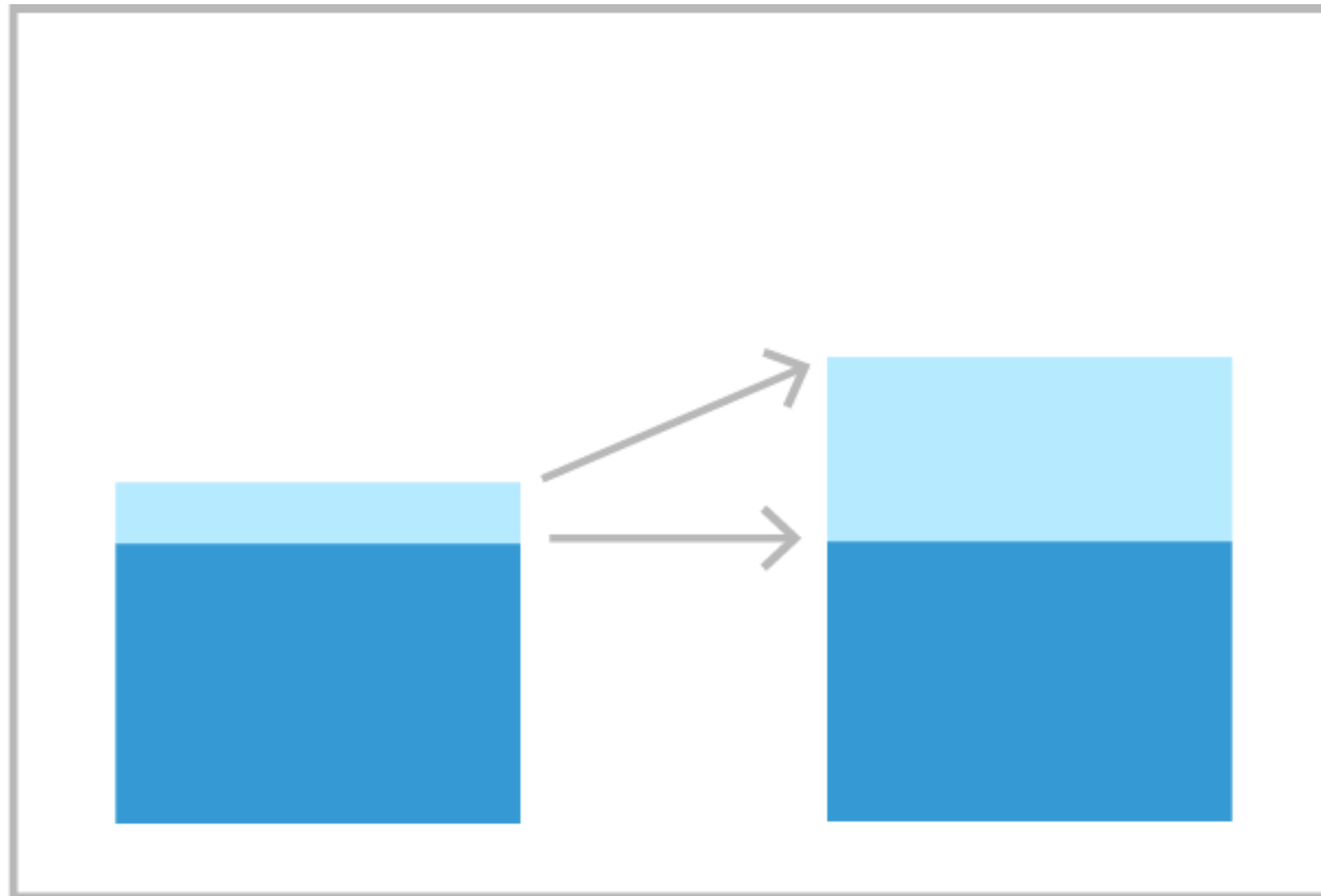


# An important use case for SDG is data amplification

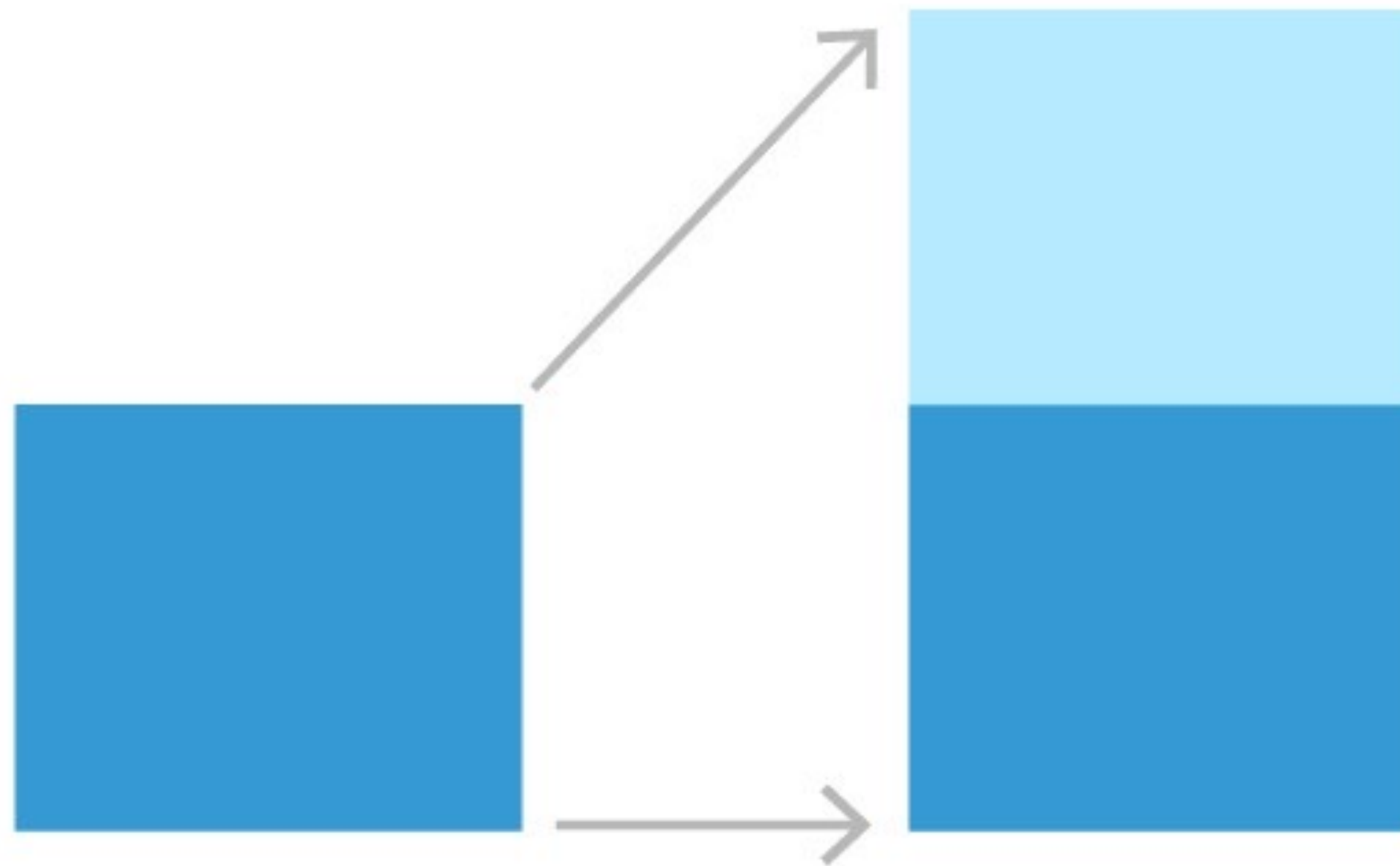




# Amplification can also focus on a specific class in the dataset

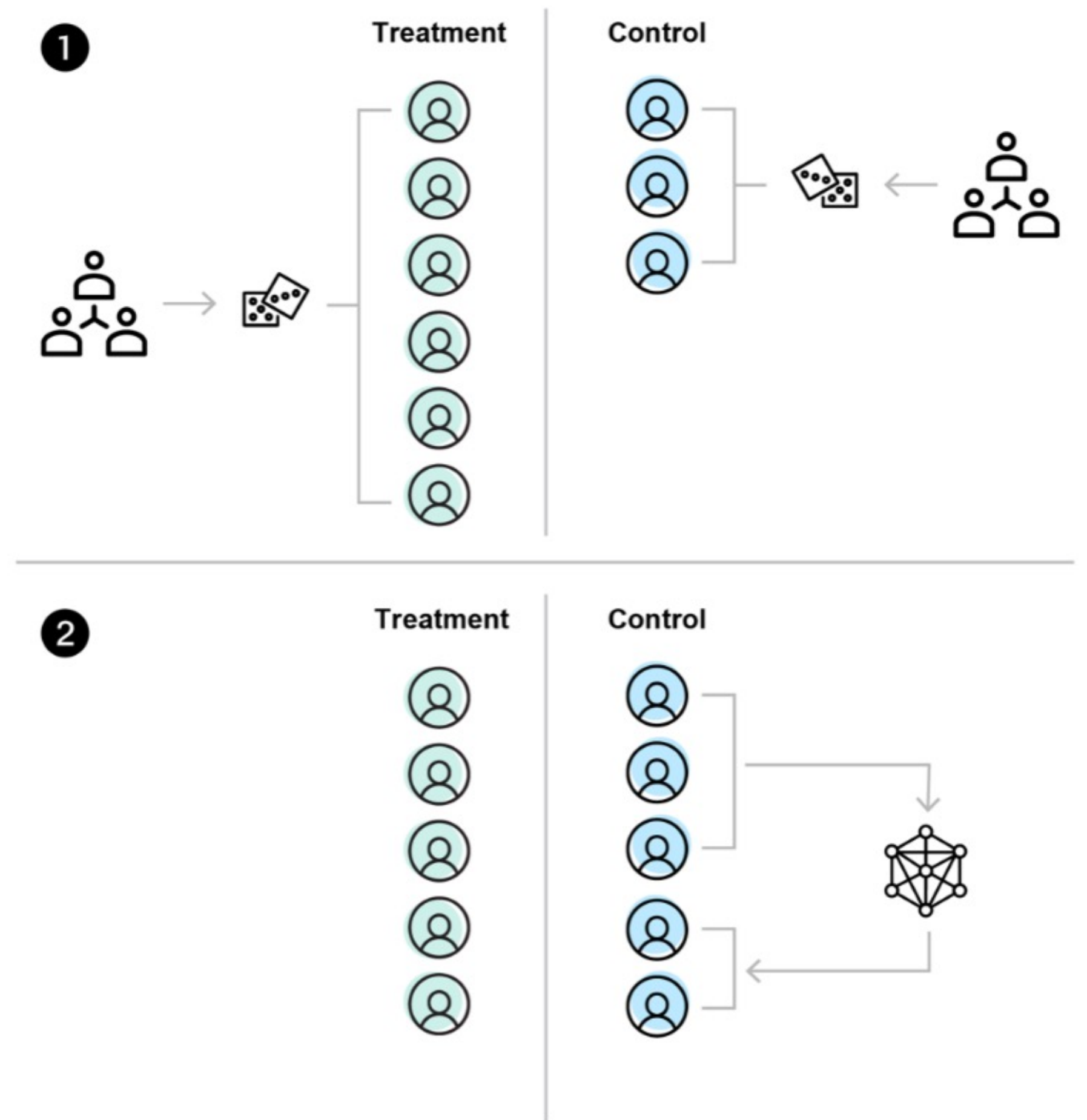


**Data augmentation is when we use SDG to simulate virtual patients to add to an existing core dataset**



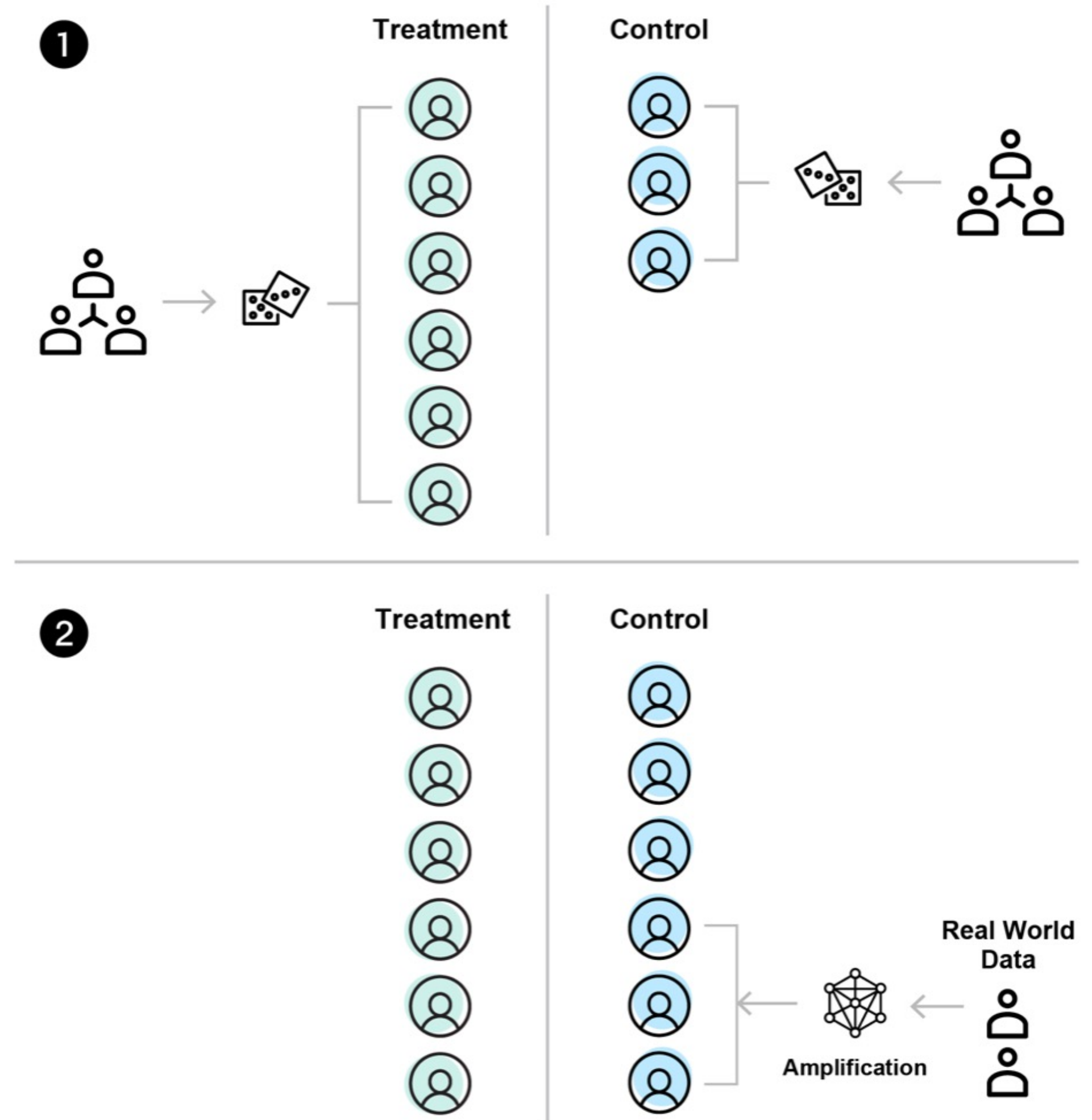
# Virtual Patients - I

Virtual patients can be simulated to reduce recruitment or to rescue studies with low recruitment or high attrition



# Virtual Patients - II

Real-world data can be amplified to create synthetic external controls, especially when there are insufficient RWD or RWD diversity





**QUESTIONS**