

# Practical Test Data Generation Adopting Deep Neural Networks

September 2021

# Agenda

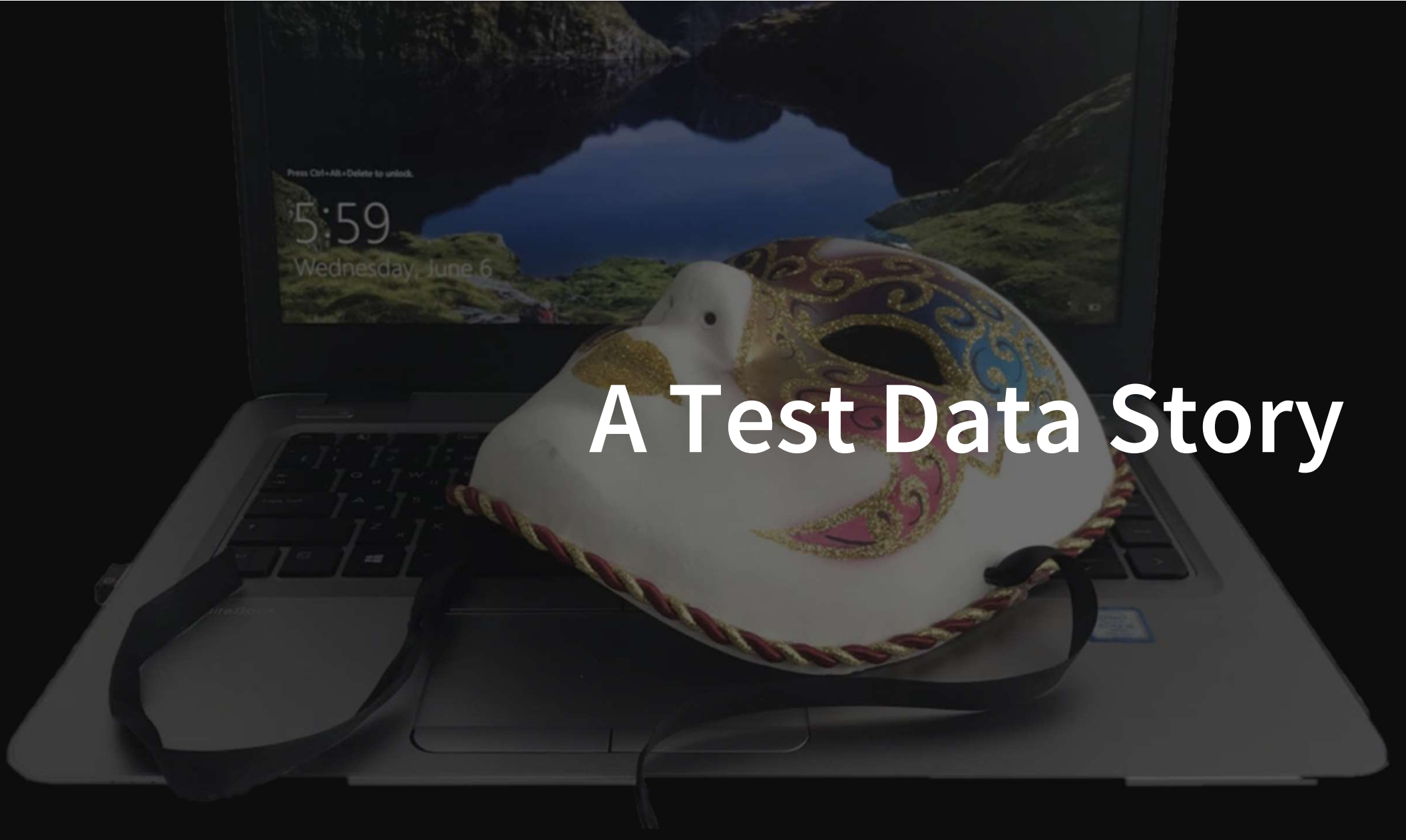
1. Intro
2. A test data story
3. From masking to machine learning
4. Beyond machine learning
5. Q&A

Press **Ctrl+Alt+Delete** to unlock.

5:59

Wednesday, June 6

# A Test Data Story



# Houston, we've had a problem!

You are going to lose access to production data in **3 months!**

1000+

Related Objects

10K+

Fields

300+

Personal and Sensitive  
Columns

---

Time to complete – **Yesterday**

# Naïve data masking



{ EI\* }
Quasi Identifiers
{ SD\* }

SSN	First Name	Last Name	Race	Date of Birth	Admission Date	Gender	ZIP	Diagnosis
123-45-6789	John	Galt	African American	1/13/1934	4/10/2018	M	80011	Testicular Cancer
772-12-4636	Horacio	Oliveira	White American	11/1/1957	4/12/2018	M	80019	Broken Arm
NULL	Natalya	Rostova	White American	10/24/1974	5/1/2018	F	80011	Pneumonia
078-05-1120	Mary	Poppins	White American	5/5/1983	5/13/2018	F	80022	Pneumonia
288-99-1020	Veena	Apsara	Asian American	6/7/1951	5/22/2018	F	80019	Broken Arm
430-09-9291	Oskar	Matzerath	White American	7/4/2018	6/3/2018	M	80011	Brain Concussion

**EI** – Explicit Identifier      **SD** – Sensitive Data

# Naïve data masking



**Quasi Identifiers**

{ EI\* }
{ SD\* }

SSN	First Name	Last Name	Race	Age, Years	Admission Date	Gender	ZIP	Diagnosis
xxx-xx-xxxx	Alfred	Uncle	African American	84	Apr, 2016	M	80000	Testicular Cancer
xxx-xx-xxxx	Benito	Gonzales	White American	61	Apr, 2016	M	80000	Broken Arm
xxx-xx-xxxx	Nadezhda	Zhdanova	White American	43	May, 2016	F	80000	Pneumonia
xxx-xx-xxxx	Deborah	Jonze	White American	35	May, 2016	F	80000	Pneumonia
xxx-xx-xxxx	Mary	Pereira	Asian American	67	May, 2016	F	80000	Broken Arm
xxx-xx-xxxx	Alex	Heilmann	White American	0	Jun, 2016	M	80000	Brain Concussion

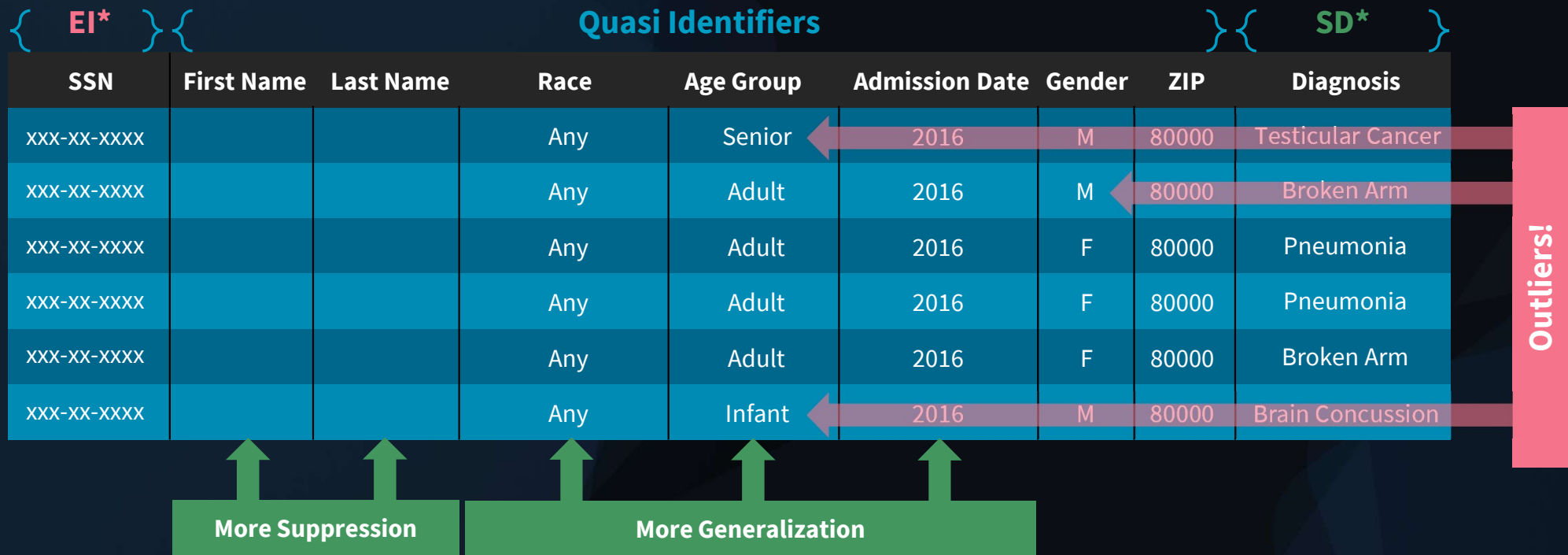
Suppression

Random Substitution

Generalization

**EI** – Explicit Identifier    **SD** – Sensitive Data

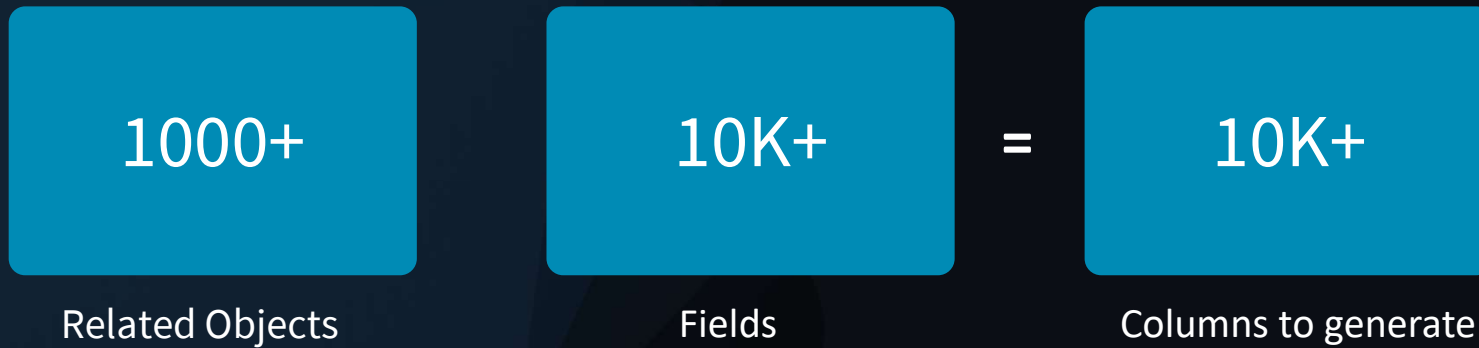
# Naïve data masking



**EI** – Explicit Identifier    **SD** – Sensitive Data

# Rules-based generation

Let get the data generated!



---

Time to complete – **Never**



# Age of Machine Learning and Artificial Intelligence



More variability is better



Not-OK data for better test coverage



Get more data than I already have



Data for testing ML/AI projects

# **Any sufficiently advanced technology is indistinguishable from magic**

Arthur C. Clarke

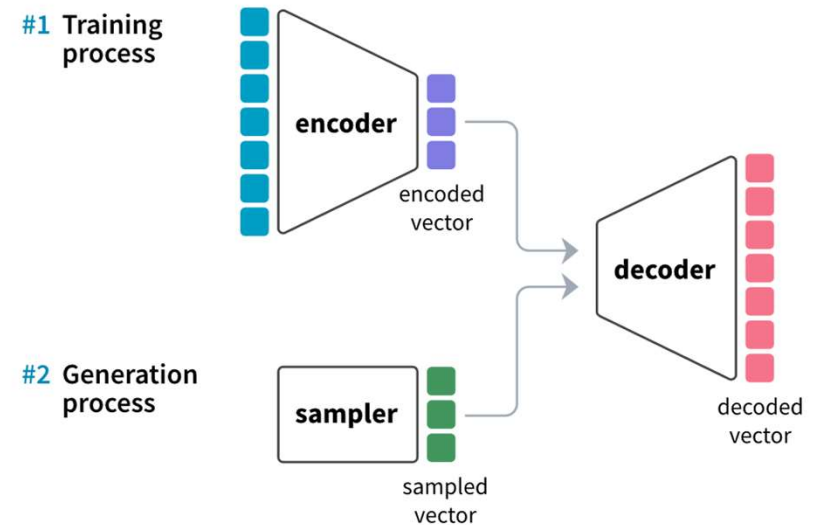


# From Masking to Machine Learning

# ML-based synthetic data generation

1. Neural networks: Variational Autoencoder (VAE), Generative Adversarial Network (GAN) and others
2. Automated process: computational resources, not people.
3. Different data types supported (categories, numeric values, free text).
4. Utility or privacy?
5. Customization is possible.
6. Various use cases:
  - generation of test data
  - data for analysis
  - data for ML engineers
  - ...

## Variational autoencoder scheme



# Original vs. synthetic

**ORIGINAL DATA**

**DATA GENERATED BY NEURAL NETWORK**

Attribute	Value	Value
Age	39	58
Workclass	State-gov	Self-emp-not-inc
Final weight	77516	134436
Education	Bachelors	Bachelors
Marital status	Never-married	Married-civ-spouse
Occupation	Adm-clerical	Craft-repair
Relationship	Not-in-family	Husband
Race	White	White
Sex	Male	Male
Hours per week	40	34
Native country	United-States	United-States
Capital	2174	3277
Income	<=50K	>50K

# Data patterns preservation

## UUID

Original	Synthetic
fe52c68a-48ec-4837-a24c-4c2bd191431f	ae92a2d1f-d233-45bb-8415-dc7b84b72477
d19953f6-1f53-4b54-8fa4-d7374ca0607b	ea4c30c0d-dc4d-45a2-ac29-3cd8877516dd
3aea0026-faae-4845-8708-9bf323ff8a77	0-7114383-fe4c-47e5-9079-33463a4b9c16

## Address

Original	Synthetic
1108 ROSS CLARK CIRCLE	4710 M EOASHA BS E ORT
145 NEWCOMB AVENUE	33 OXNITTHOH2 X1VT
2540 EAST ST	150 1BETM EY STA H TNNEBE

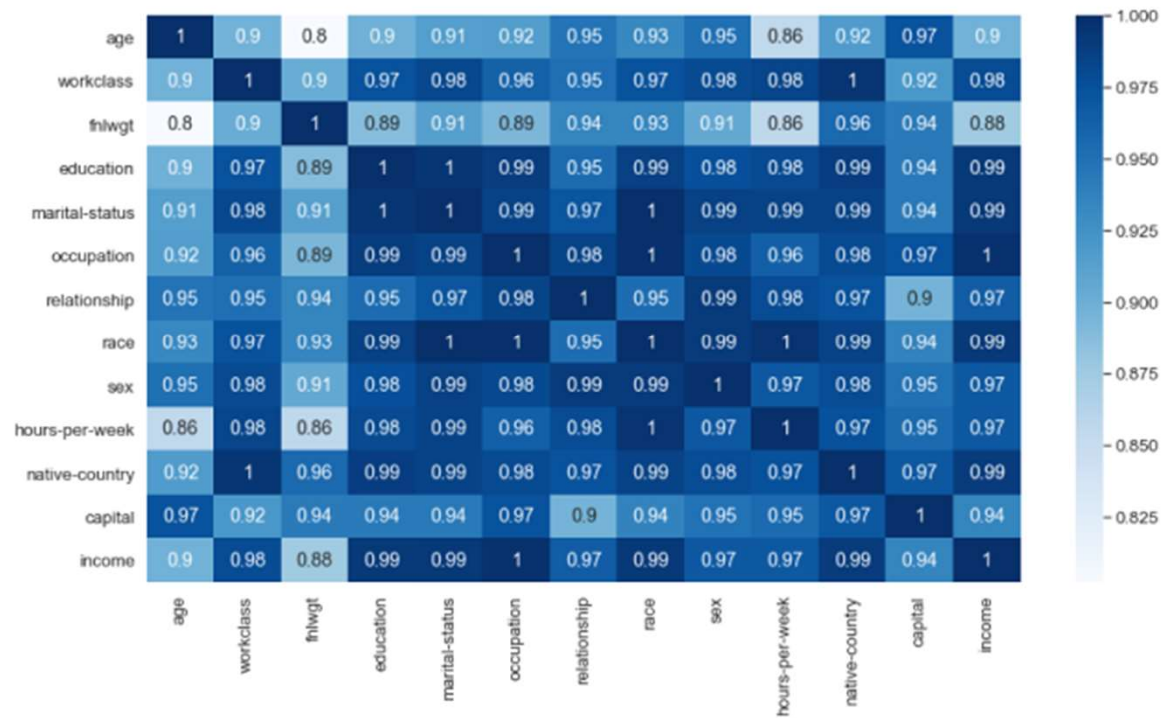
## Serial number

Original	Synthetic
S000028330	S000023777
S000026892	S000000850
S000005529	S000043975

# Is new synthetic data useful?

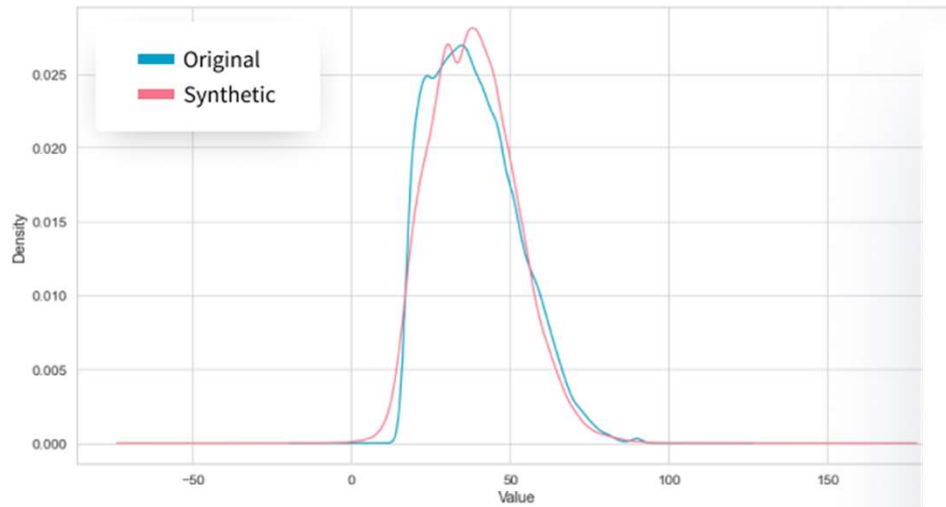
Median accuracy: 0.969

Accuracy: original ~ synthetic

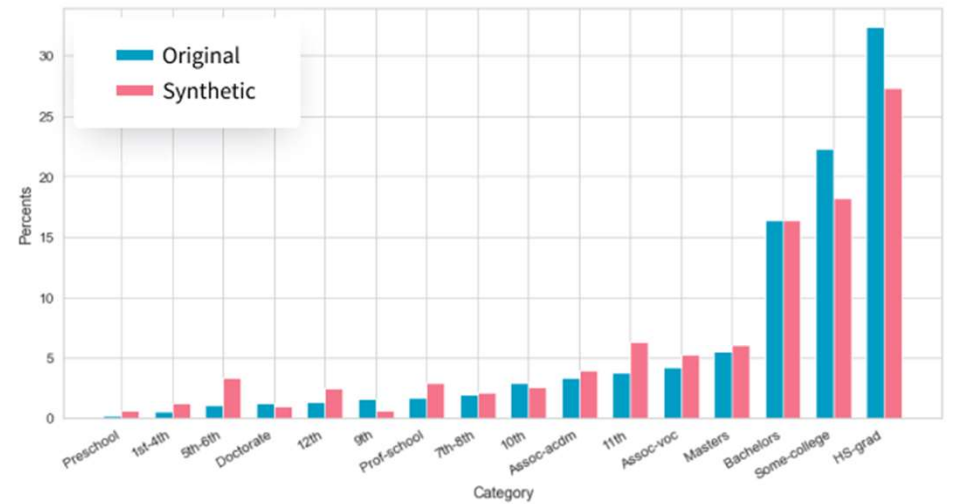


# Univariate distributions

## Age



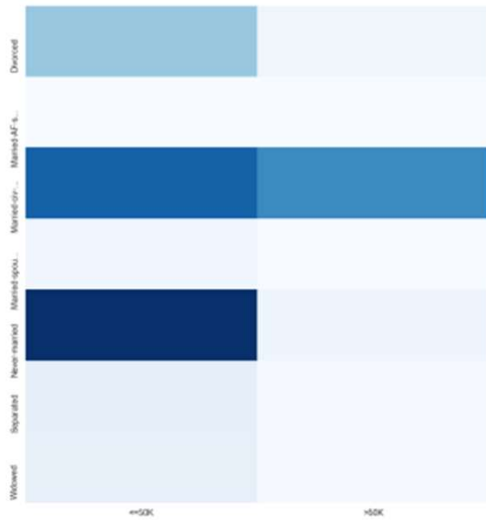
## Education



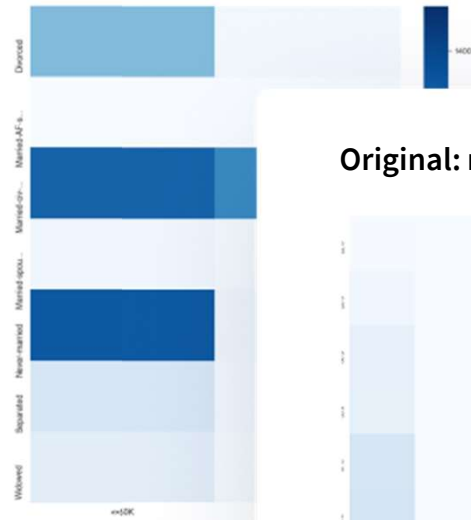


# Bivariate distributions

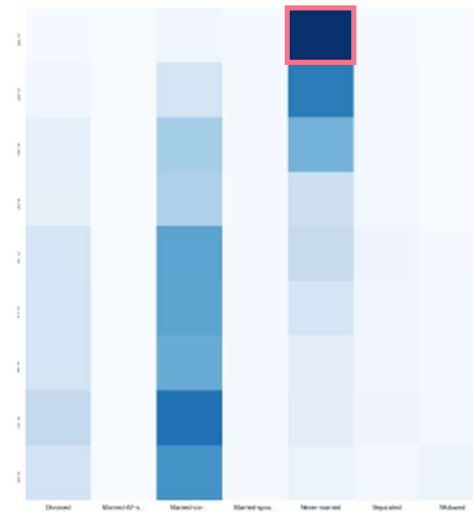
Original: age vs marital-status



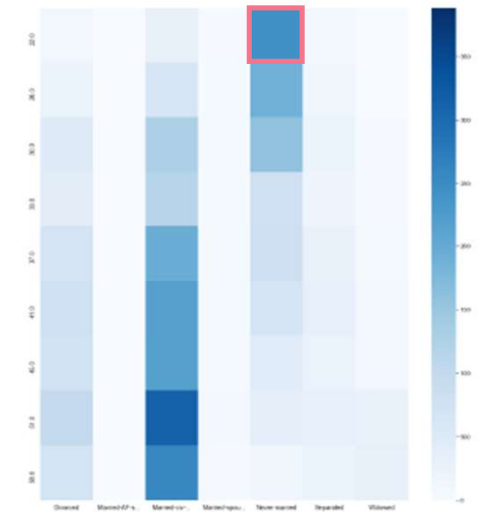
Synthetic: age vs marital-status



Original: marital-status vs income

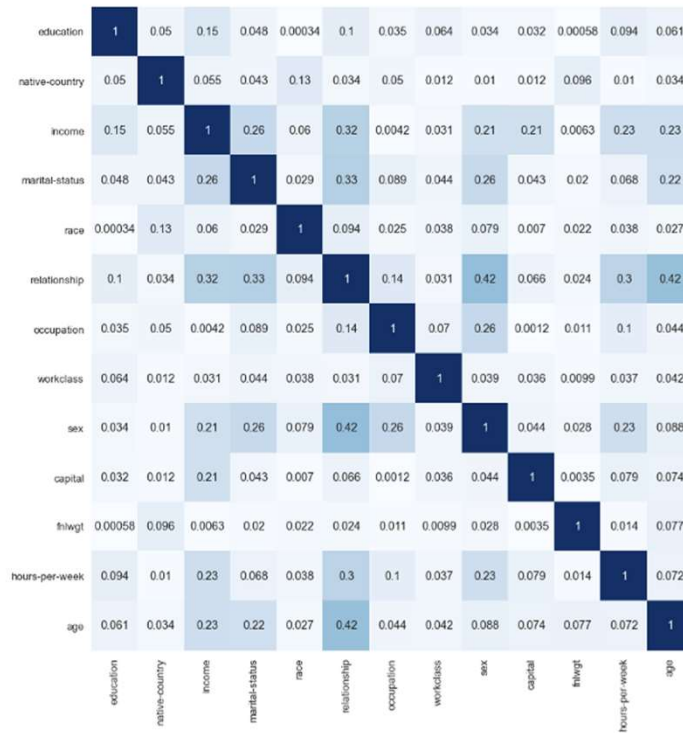


Synthetic: marital-status vs income

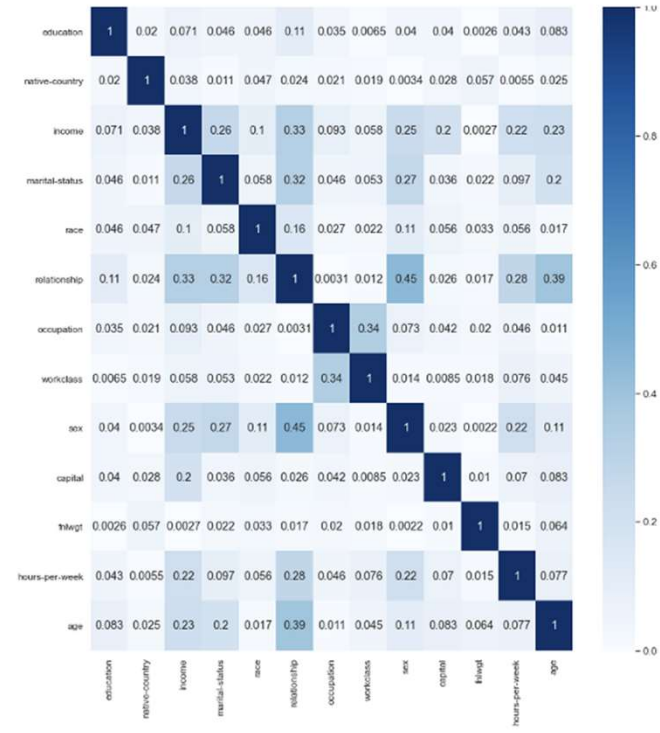


# Correlations between columns

Original



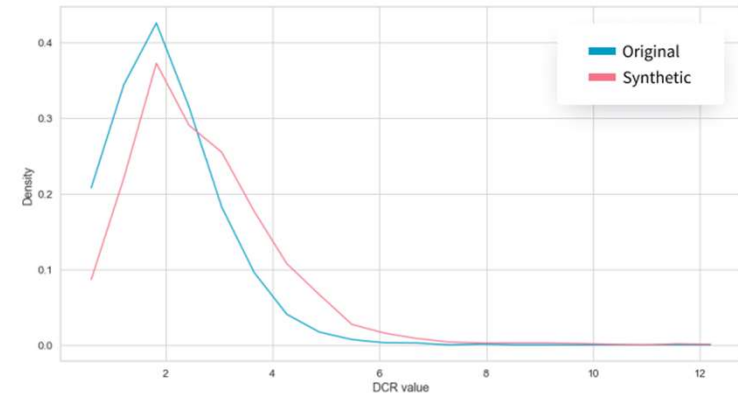
Synthetic



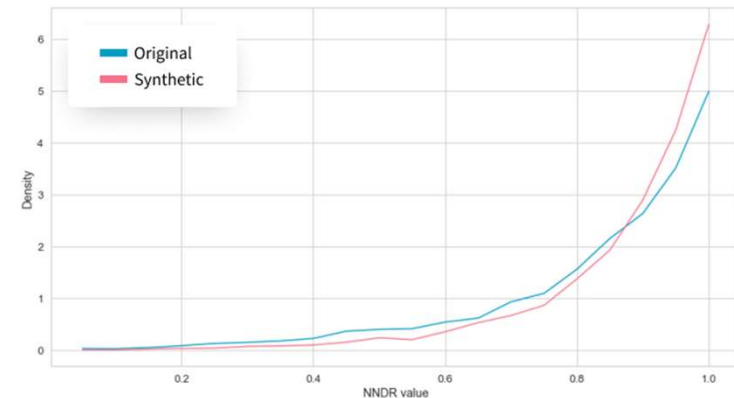
# Is new synthetic data private?

1. Distance-based metrics (comparing original hold-out and generated data):
  - Full matches
  - Distance to closest record (DCR)
  - Nearest neighbor distance ratio (NNDR)
2. ML-based (training an attacker model)

Passed DCR test



Passed NNDR test



# Need more privacy?

**Differential privacy** – a methodology of making computations private with quantifiable and provable worst case scenario leakage of the original data.

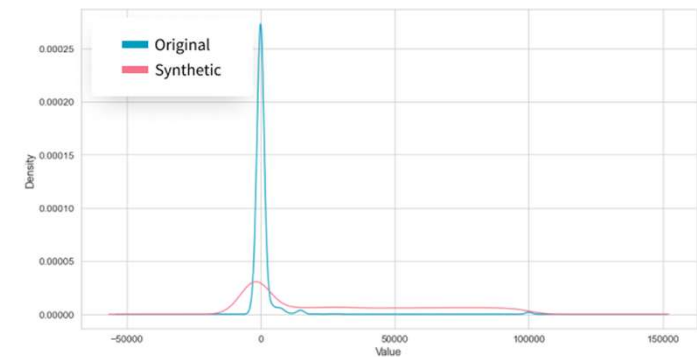
Variational Autoencoder with differentially private training process: **DP-SGD**

Privacy level can be controlled with parameter  $\epsilon$  (privacy budget).

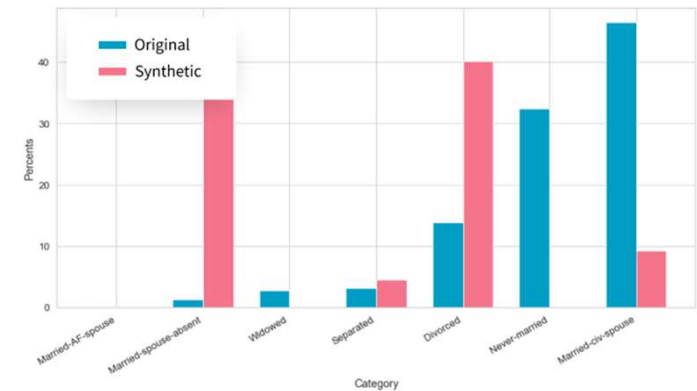


**Capital**

$\epsilon = 1.4$



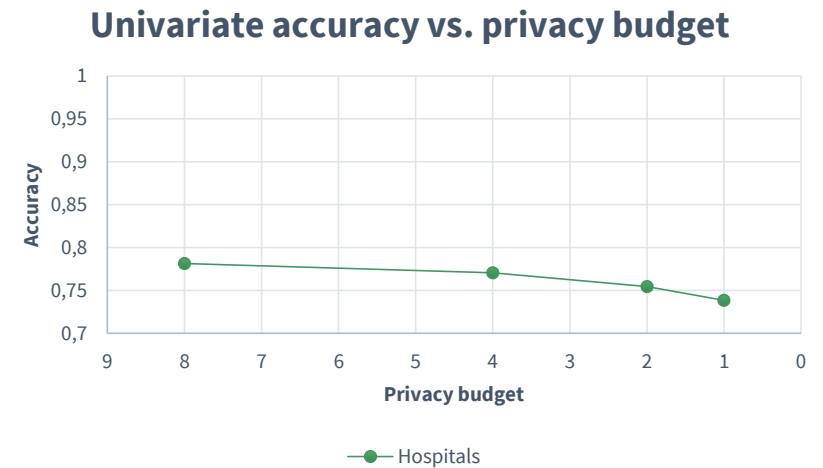
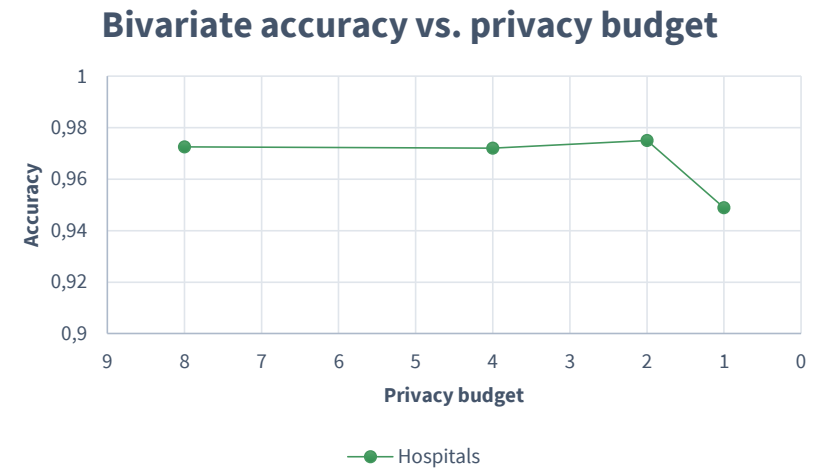
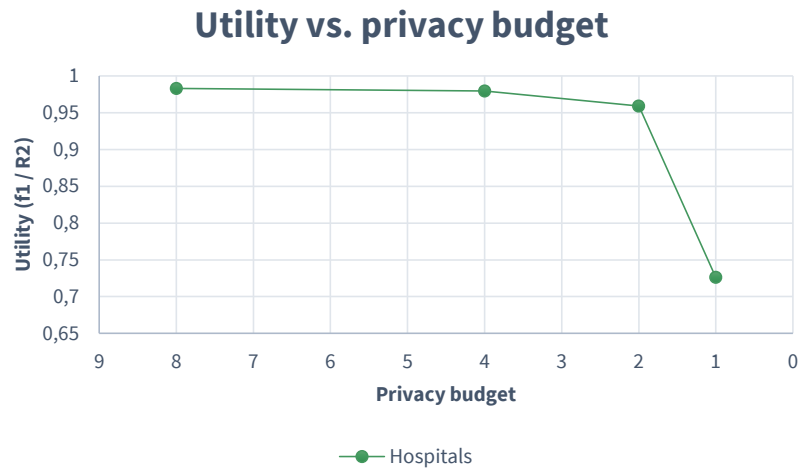
**Marital status**



# What is the utility/privacy exchange rate?

## The experiments show that

- Considerable data utility drop happens for privacy budgets less than 2
- In some cases, the utility of the generated data might be higher than the original



# What do I need to get the synthetic data?

## Dataset with hospital data:

- ~50000 rows
- 5 free text columns
- 7 numeric columns
- 7 categorical columns

## Dataset with us census income data:

- ~50000 rows
- 4 numeric columns
- 9 categorical columns

**Google Cloud VM shape:** n1-standard-16 (16 vCPU, 60 GB RAM)

**Accelerator:** 1x NVIDIA TESLA P100

**Hourly rate:** \$2.27

Dataset	Epochs	Rows	Time	Minute/Epoch	Cost
Hospital: complications and deaths	200	50000	22h 1m	~6m 25s	\$50
US census income	200	50000	7h 53m	~2m 20s	\$18

# What do I need to get the synthetic data?

## Machine specifications:

- **CPU:** Intel(R) Core(TM) i5-10310U CPU @ 1.70GHz, 2208 Mhz, 4 Core(s), 8 Logical Processor(s)
- **RAM:** 32 GB

## Hospital: complications and deaths:

- 5 free text columns
- 7 numeric columns
- 7 categorical columns

Rows	Size (KB)	RAM	Time	Rows/s
1,000	320	400	1m 2s	31
10,000	3200	630	5m 12s	32
100,000	32000	2200	1h 20m	21

## US census income:

- 4 numeric columns
- 9 categorical columns

Rows	Size (KB)	RAM	Time	Rows/s
1,000	112	400	1s	1000
10,000	1097	400	1s	10000
100,000	11220	400	4s	25000

# Training of synthetic data generation models

Resources required for training **scale liner** and depend on:

- Size of the dataset – number of rows and columns
- Types of the columns: texts, categorical, numeric
- Number of epochs
- Differential privacy settings

Training specs	
Platform	AWS SageMaker
Shape	ml.g4dn.xlarge
Cost, USD/hour	0.74

<b>Number of rows</b>	10000	5000	10000	5000	500
<b>Differential privacy</b>	Y	Y	N	N	N
<b>Time, hours</b>	3	1.5	0.2	0.08	0.06
<b>Cost for 20 epochs, USD</b>	<b>2.22</b>	<b>1.11</b>	<b>0.12</b>	<b>0.03</b>	<b>0.02</b>
<b>Cost per epoch</b>	0.11	0.06	0.006	0.0015	0.001
<b>Estimated cost for 50 epochs</b>	5.5	3	0.3	0.075	0.05

## Example:

The cost of a model training on the dataset with **100K** rows and **15** columns is **2.22 \* 10 = 22 USD** in AWS SageMaker

**Note:** Prices for the regular GPU-enabled VMs are lower but require careful planning of the workload. We recommend using on-spot instances for non-time-critical calculations.

Data set	
<a href="#">Churn_modelling.csv</a>	
Columns	14
Text	1
Numeric	6
Categorical	7



# Synthetic data generation estimations

Resources required for generation **scale liner** and depend on:

- Size of the dataset – number of rows and columns
- Types of the columns: texts, categorical, numeric

Generation specs	
Platform	AWS SageMaker
Shape	ml.m5.large
Cost, USD/hour	0.12

<b>Number of Rows</b>	5000	2000	500
<b>Time, hours</b>	0.1	0.08	0.08
<b>Cost, USD</b>	<b>0.012</b>	<b>0.0096</b>	<b>0.0096</b>

## Example:

The cost of a dataset generation with **100K** rows and **15** columns is  **$0.012 * 20 = 0.24$**  **USD** in AWS SageMaker

**Note:** Data generation does not require GPU-enabled VM and can be done on most types of machines, including workstations.

Data set	
<a href="#">Churn_modelling.csv</a>	
Columns	14
Text	1
Numeric	6
Categorical	7



# Beyond ML

## Trained models as a data source

<b>Train</b>	Train one model per table
<b>Relate</b>	Pass information on relationships between tables to the models
<b>Preserve</b>	Preserve data types of the attributes
<b>Manage</b>	Manage set of trained models as a single data source
<b>Evaluate</b>	Evaluate models performance after training
<b>Control</b>	Control the training cost
<b>Generate</b>	Generate data and load it into various data stores

# Create new pipeline

From the groups of related tables, select one table as a starting point for data transfer. The graph displays the selected table and its parents and children. Choose the required incoming and outgoing link depth to regulate how many tables will be included

ISOLATED GRAPHS (243) Search by table name

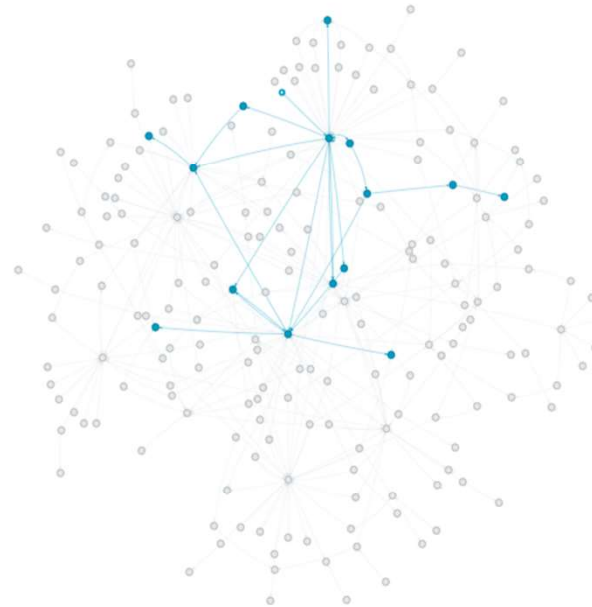
Manage manual relationships

Incoming depth:  Outgoing depth:

Name	
Group 1	194
bgb_aanleveringen	B
bgb_afdelingen	
bgb_agp_classificaties	
bgb_apparatuur	
bgb_areas	
bgb_artikelgroep_landen	
bgb_artikelgroep_statistieken	
bgb_artikelgroep_talen	
bgb_artikelgroep_xrefs	
bgb_artikelgroepen	
bgb_artikelomschrijvingen	
bgb_ass_instellingen	
bgb_assortimentsgroepen	
bgb_assortimentgroep_artgroups	
bgb_assortimentgroep_shelve_afs	
bgb_assortimentgroups	

Base table: Group 1 • bgb\_aanleveringen

Hold Sh



## Visualize

— Predefined relationships —

Cancel Save as draft < Back

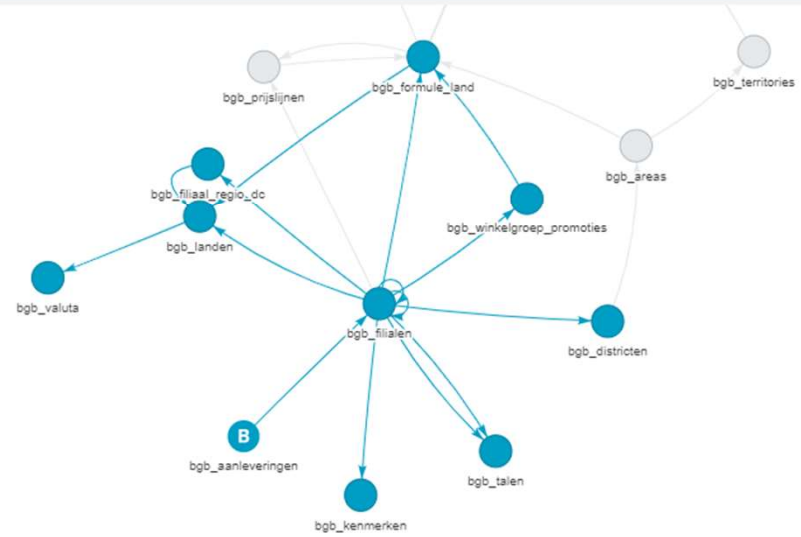
# Create new pipeline

✓ Set up basics > ✓ Design a subset > 3 Refine metadata

Now you can apply filters, exclude unnecessary tables or columns from your dataset.

TABLES (10/16) Search by table name

<input type="checkbox"/> Name	<input type="checkbox"/> Hide unchecked
<input checked="" type="checkbox"/> bgb_aanleveringen	
<input type="checkbox"/> bgb_areas	
<input checked="" type="checkbox"/> bgb_districten	
<input checked="" type="checkbox"/> bgb_filiaal_regio_dc	
<input checked="" type="checkbox"/> bgb_filialen	
<input type="checkbox"/> bgb_formula_country_groups	
<input checked="" type="checkbox"/> bgb_formule_land	
<input type="checkbox"/> bgb_formules	
<input checked="" type="checkbox"/> bgb_kenmerken	
<input checked="" type="checkbox"/> bgb_landen	
<input type="checkbox"/> bgb_prijlijnen	
<input checked="" type="checkbox"/> bgb_talen	
<input type="checkbox"/> bgb_territories	
<input checked="" type="checkbox"/> bgb_valuta	



COLUMNS (0/0) RELATIONSHIPS (0/1)

Search by column name

## Subset

Column name

Connections

Type

Filters (0)

Cancel

Save as draft

< Back

# Customer Data Model

Trained

Train

Edit model



## DESCRIPTION

In this guide, we'll make some calls to the GitHub Enterprise Server Search API, and iterate over the results using pagination. You can find the complete source code for this project in the platform-samples repository

<b>TYPE</b> Data subsetting	<b>FROM</b> Source Data Base Name	<b>WHERE</b> Cluster name	<b>WHAT</b> Order_Items, 45 tables of 99
--------------------------------	--------------------------------------	------------------------------	---

Training activity Metadata

Display 20 on page

Showing 1-20 of 57 1 2 3 >

Start time	Finish time	Duration	Started by	Status	Actions							
2021-02-09 13:15	--	--	Superadmin	In progress	Download	Metrics report	Restore					
<div style="border: 1px solid #007bff; padding: 5px;"><p>Data subsetting process started...</p><table border="1"><tr><td><b>TYPE</b> Data subsetting</td><td><b>FROM</b> SourceDataBaseName</td><td><b>TO</b> Old Target Name</td><td><b>WHERE</b> Cluster name</td><td><b>WHAT</b> Order_Items, 45 tables of 99</td></tr></table><p><b>DESCRIPTION</b> In this guide, we'll make some calls to the GitHub Enterprise Server Search API, and iterate over the results using pagination. You can find the complete source code for this project</p></div>								<b>TYPE</b> Data subsetting	<b>FROM</b> SourceDataBaseName	<b>TO</b> Old Target Name	<b>WHERE</b> Cluster name	<b>WHAT</b> Order_Items, 45 tables of 99
<b>TYPE</b> Data subsetting	<b>FROM</b> SourceDataBaseName	<b>TO</b> Old Target Name	<b>WHERE</b> Cluster name	<b>WHAT</b> Order_Items, 45 tables of 99								
2021-02-09 13:15	2021-02-04 11:59	00:03:56	Superadmin	Done	Download	Metrics report	Restore					
2021-02-04 11:55	2021-02-04 11:59	00:03:56	Superadmin	Done	Download	Metrics report	Restore					
2021-02-09 13:00	2021-02-09 15:03	00:04:01	Superadmin	Failed	Download	Metrics report	Restore					

Manage

## Likely limitations

---

Dependencies between relational entities might be broken

---

Too much training data increases requirements to the hardware resources

---

Too little training data affects accuracy of the generated data

---

The model maintains correlations within datasets but not with the outside world objects

---

Machine readable formats such as JSON and XML

# Thank you!



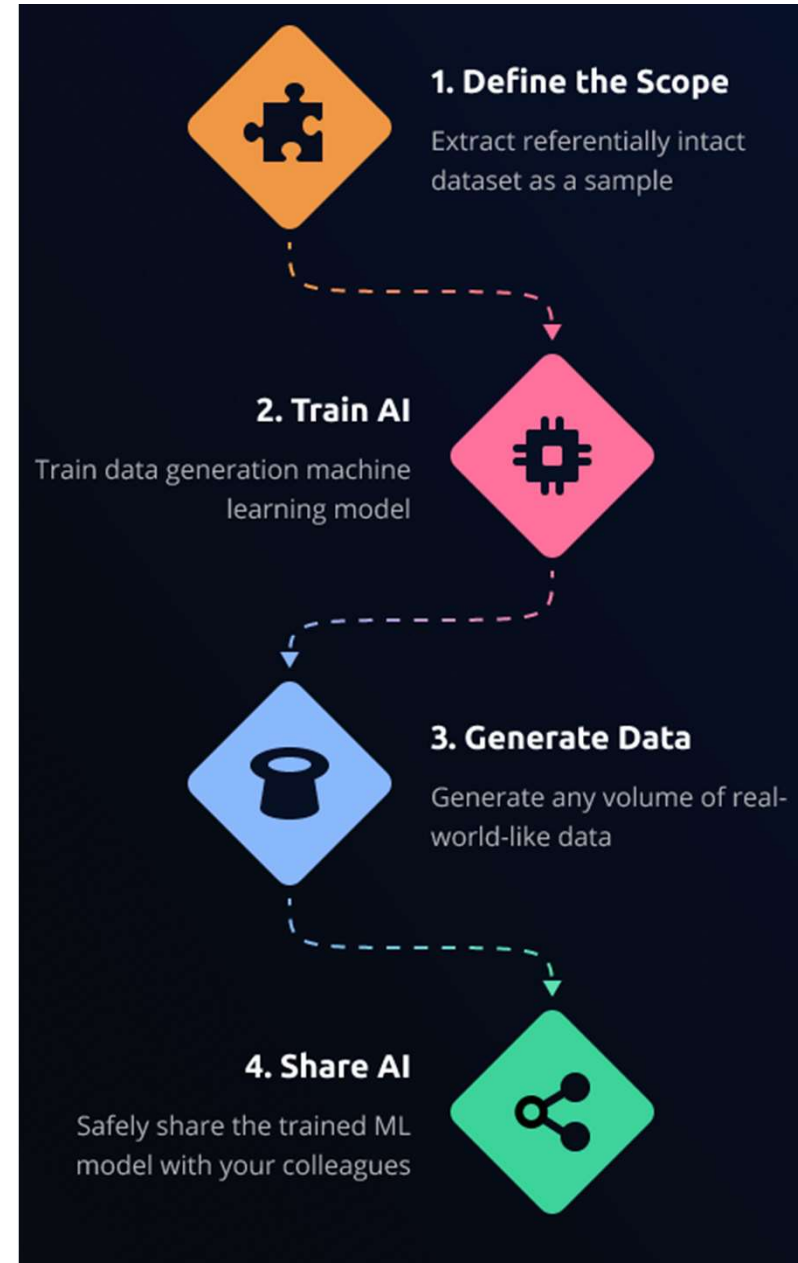
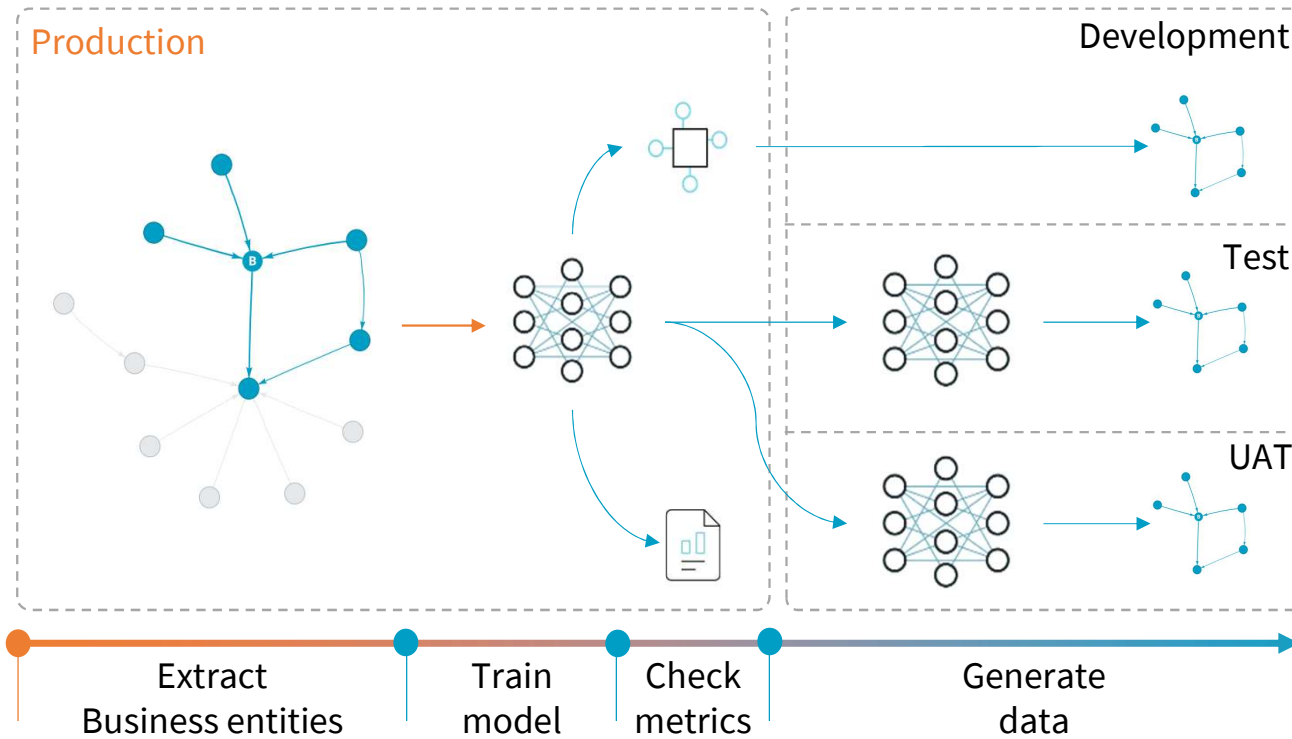
For more information, please contact:  
[supportepmc-tdm@epam.com](mailto:supportepmc-tdm@epam.com)



# Appendix

# High level flow

**PRODUCTION DATA STAYS WITHIN SECURED PERIMETER**



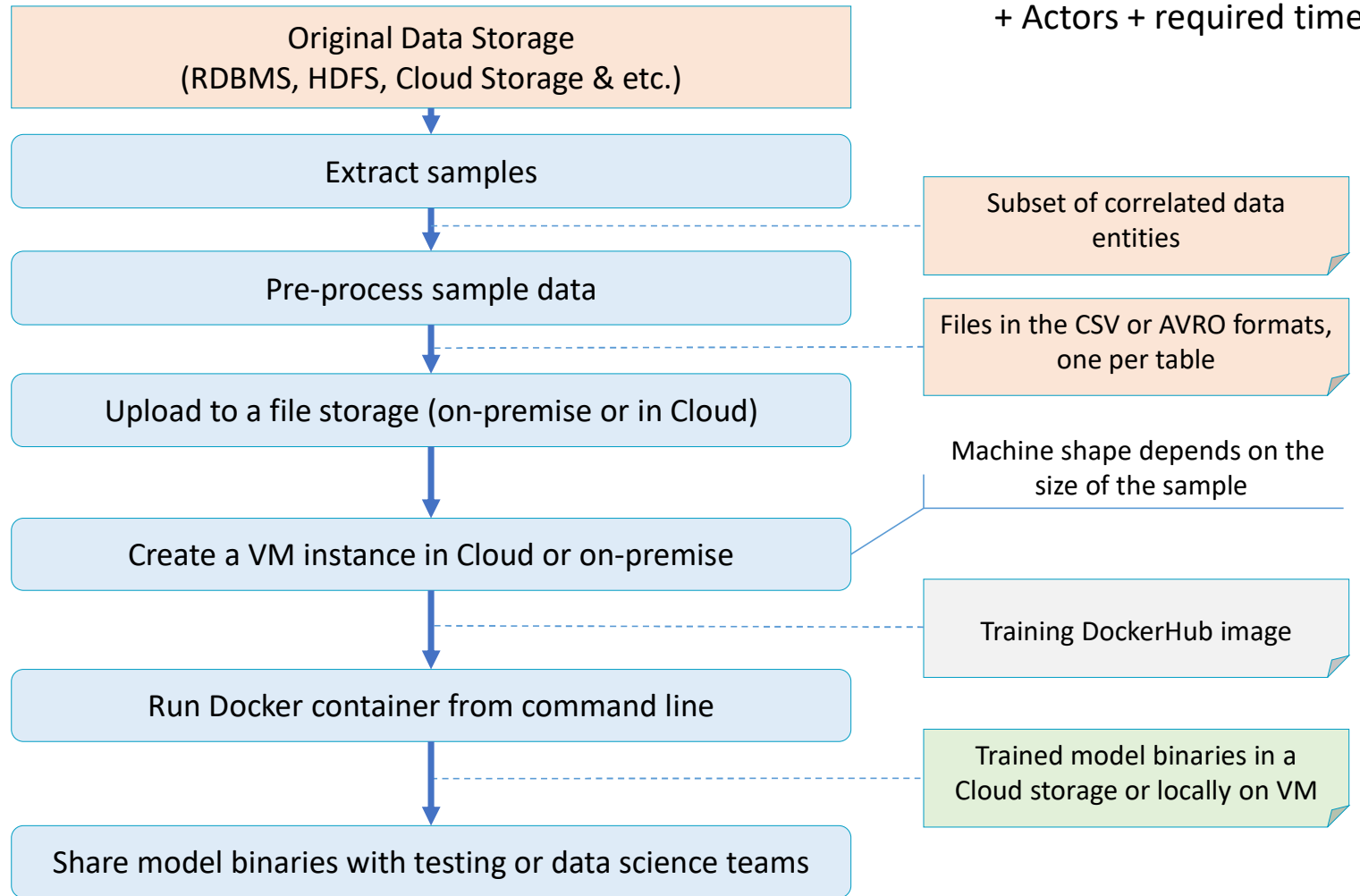
# Reference implementation: Training under your control

+ Actors + required time

**Preparation**

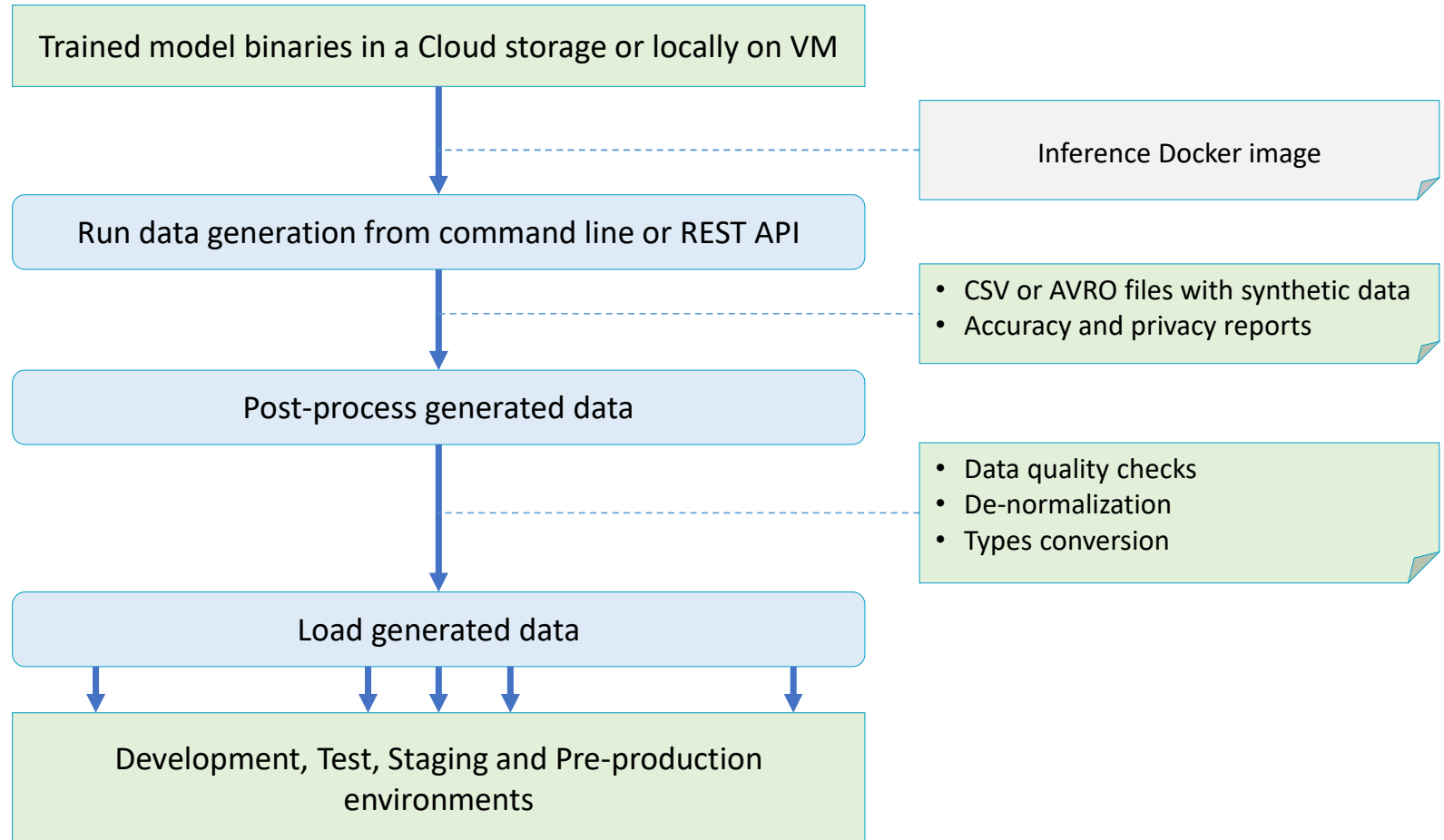


**Model training**



# Reference implementation: Synthetic data generation

## Generation



## Compare with the “Classic” flow

### CLASSIC

#### Pros

- Precision and control over the data classification and mutation
- High data utility preservation (masking)
- Flexible rules and techniques
- High performance and low resource consumption
- Consistency and repeatability

#### Cons

- Labor-intensive steps
- Risks associated with a human error
- Requires training of personnel to use the tool efficiently

### ML-DRIVEN

Experimental

#### Pros

- Short test data delivery cycle
- No or minimal human involvement
- Preserves hidden patterns in data and suitable for ML tasks
- Allows moving models between security perimeters rather than data
- Comprehensive metrics
- Increases test data variability and injects boundary cases

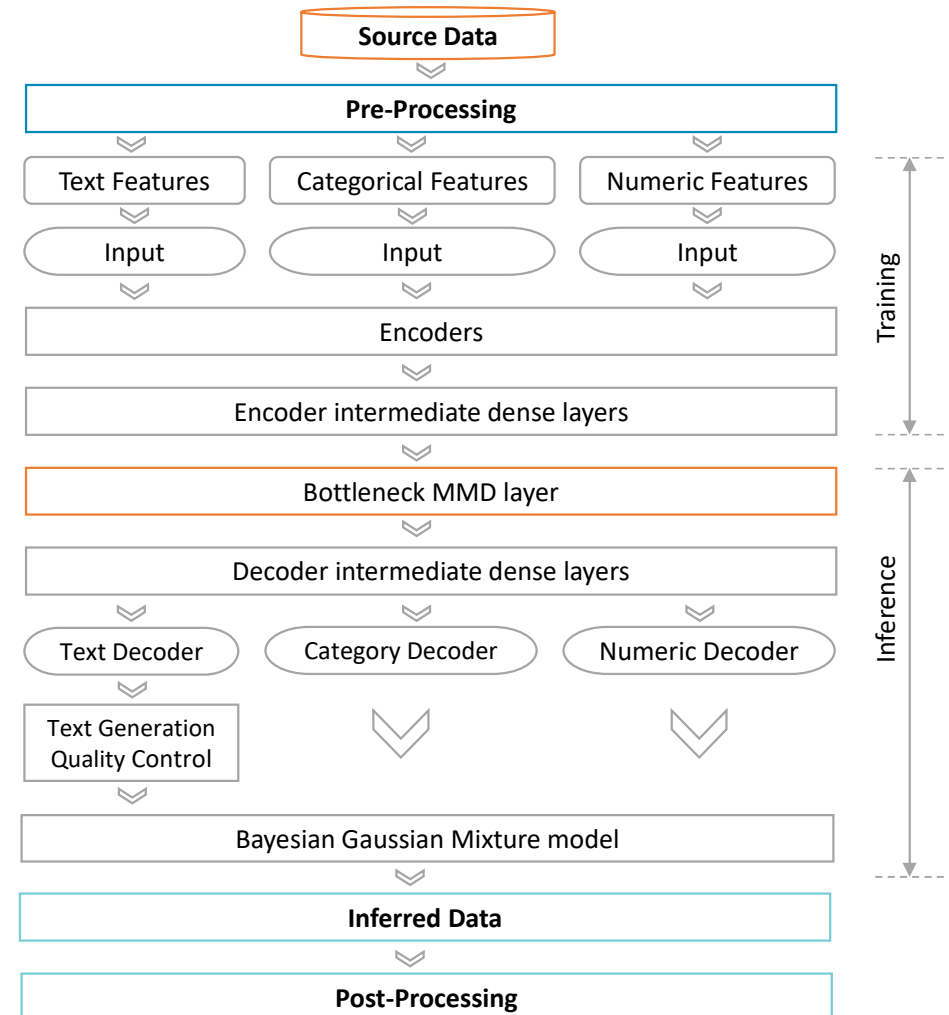
#### Cons

- High resource consumption for training and inference
- No consistency or repeatability
- Limited support of the relationships between datasets

## Internal Model Architecture

In the core of the data generation module lies a customized deep learning model, based on the Info VAE (Variational Autoencoder) with the following modifications:

- MMD bottleneck layer is added to regularize the latent space
- Bayesian gaussian mixture model (BGMM) on the latent space identifies latent Gaussian regimes
- Unsupervised datatypes detection is employed for preprocessing pipeline
- Model transformers are deployed to generate complex text when needed
- Differential Privacy (DP) enabled as part of the model training process



## Samples – Investment Series

### ORIGINAL

Attribute	Value
Reporting File Number	811-00582
CIK	44402
Name of Registrant	NEUBERGER BERMAN EQUITY FUNDS
Org Type	30
Series ID	S000007847
Series Name	Neuberger Berman Socially Responsive Fund
Address_1	605 THIRD AVENUE
City	NEW YORK
State	NY
Reporting File Number	811-21714
CIK	1317146
Name of Registrant	MML Series Investment Fund II
Org Type	30
Series ID	S000028330
Series Name	MML Short-Duration Bond Fund
Address_1	1295 STATE STREET
City	SPRINGFIELD
State	MA

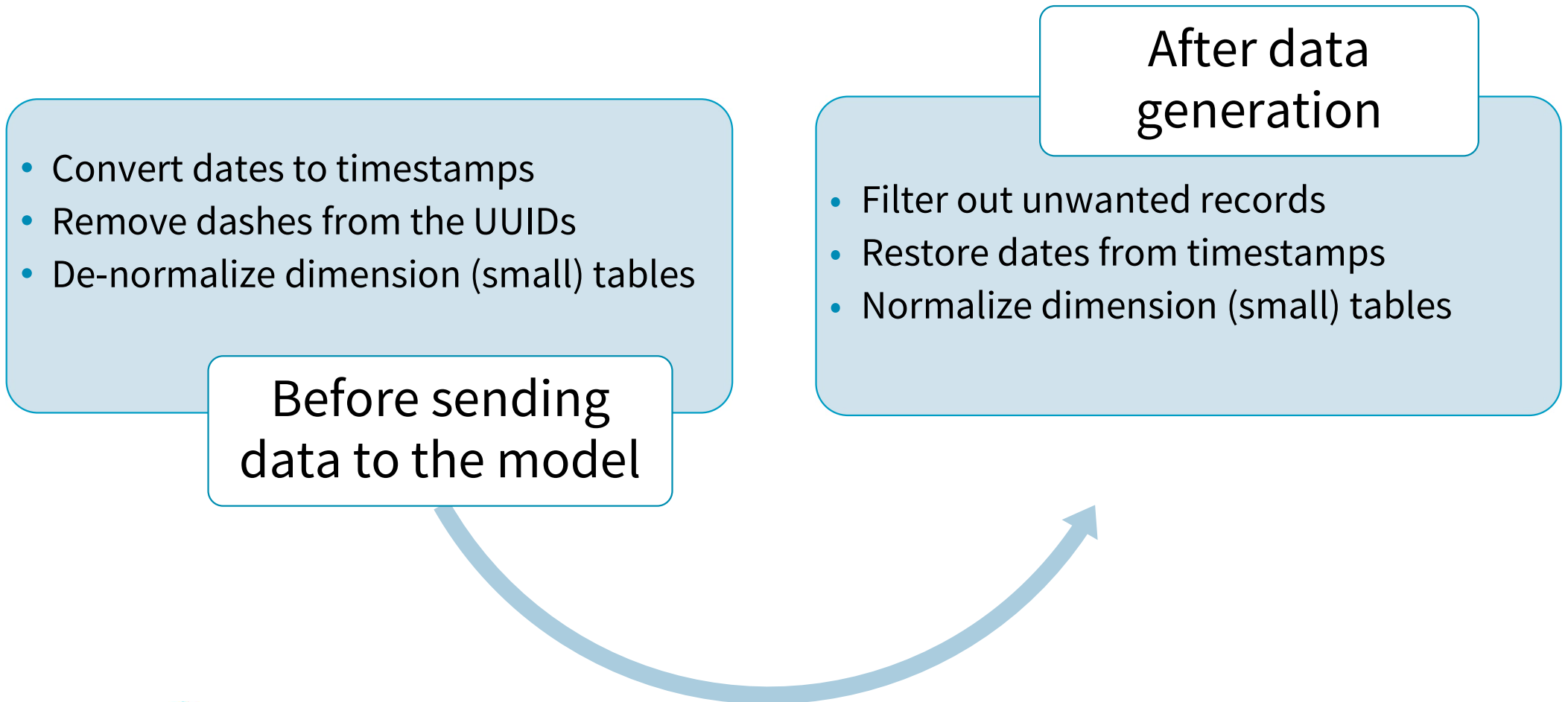
### GENERATED

Attribute	Value
Reporting File Number	811-00314
CIK	204804
Name of Registrant	TRRESERWESREE M F FN NUNNS
Org Type	30
Series ID	S000023777
Series Name	GranmaaeEiHT MMh ti roaTemnad ectFe Ftnd
Address_1	1523 SOOHH TUS RT
City	DRNTENONI LLLLLGE
State	CO
Reporting File Number	811-22427
CIK	1391901
Name of Registrant	ToS e slA dsa r iusr
Org Type	30
Series ID	S000038281
Series Name	B..Rn nRIrriSplo On2pdeet dai0ul uF
Address_1	101 EAUH TTON SSRRET
City	BALTIMORE
State	MD

Inferred Format

Preserved Text Structure

## Pre- and post- processing as the way to improve the accuracy





## References on privacy metrics

- «Fidelity and Privacy of Synthetic Medical Data (Review of Methods and Experimental Results)», Ofer Mendelevitch, Michael D. Lesh SM MD FACC, июнь 2021  
<https://arxiv.org/ftp/arxiv/papers/2101/2101.08658.pdf>
- «Bootstrap confidence intervals», Jeremy Orloff and Jonathan Bloom  
[https://ocw.mit.edu/courses/mathematics/18-05-introduction-to-probability-and-statistics-spring-2014/readings/MIT18\\_05S14\\_Reading24.pdf](https://ocw.mit.edu/courses/mathematics/18-05-introduction-to-probability-and-statistics-spring-2014/readings/MIT18_05S14_Reading24.pdf)
- Short overview of metrics above  
<https://mostly.ai/2020/11/04/truly-anonymous-synthetic-data-legal-definitions-part-ii/>
- Implementation of quality, accuracy, and privacy metrics for the ML data generation models  
[https://sdv.dev/SDV/user\\_guides/evaluation/index.html](https://sdv.dev/SDV/user_guides/evaluation/index.html)  
<https://github.com/sdv-dev/SDMetrics>

## Some examples

Original table:

Age	Workclass	fnlwght	Education	Marital status	Occupation	Relationship	Race	Sex	Hours per week	Native country	Capital	Income
39	State-gov	77516	Bachelors	Never-married	Adm-clerical	Not-in-family	White	Male	40	United-States	2174	<=50K
50	Self-emp-not-inc	83311	Bachelors	Married-civ-spouse	Exec-managerial	Husband	White	Male	13	United-States	0	<=50K
28	Private	338409	Bachelors	Married-civ-spouse	Prof-specialty	Wife	Black	Female	40	Cuba	0	<=50K

Table generated by neural network:

Age	Workclass	Fnlwght	Education	Marital status	Occupation	Relationship	Race	Sex	Hours per week	Native country	Capital	Income
32	Private	249175	Some-college	Never-married	Other-service	Not-in-family	White	Female	40	United-States	40	<=50K
50	Private	218399	Bachelors	Married-civ-spouse	Exec-managerial	Husband	White	Male	48	United-States	-54	>50K
25	Local-gov	114298	Bachelors	Divorced	Prof-specialty	Unmarried	White	Female	34	Canada	40	<=50K