

*Antreas Kapenakakis @ Aalborg University, antreas@cs.aau.dk*

**Data privacy** as a topic has entered the mainstream conscious. Years of negligent data storage, leading to breaches, and **inappropriate sharing** have made the public reluctant to share their data. In addition, this caused legislators to enact **restrictions** in the **use and handling of data**.

## Data Release Types

Data releases can be broadly split into 2: **datasets** and **the results of analyses** of the data.

Publishing the **results of an analysis** (models, statistics, census data) lowers the possible exposure of the participants but **limits the usability** of the data to the **questions answered by the release**. However, since the information output is limited it is **easier to protect** the data using a tool such as **Differential Privacy**.

A more general **data release**, such as a **dataset**, allows for answering any arbitrary question. This larger scope raises questions about **which information** should be kept in the data, and how it should be **evaluated in terms of privacy and accuracy**.

## Ensuring Privacy and Accuracy

In order to ensure that a **data release** meets **accuracy** and **privacy requirements**, two approaches are used.

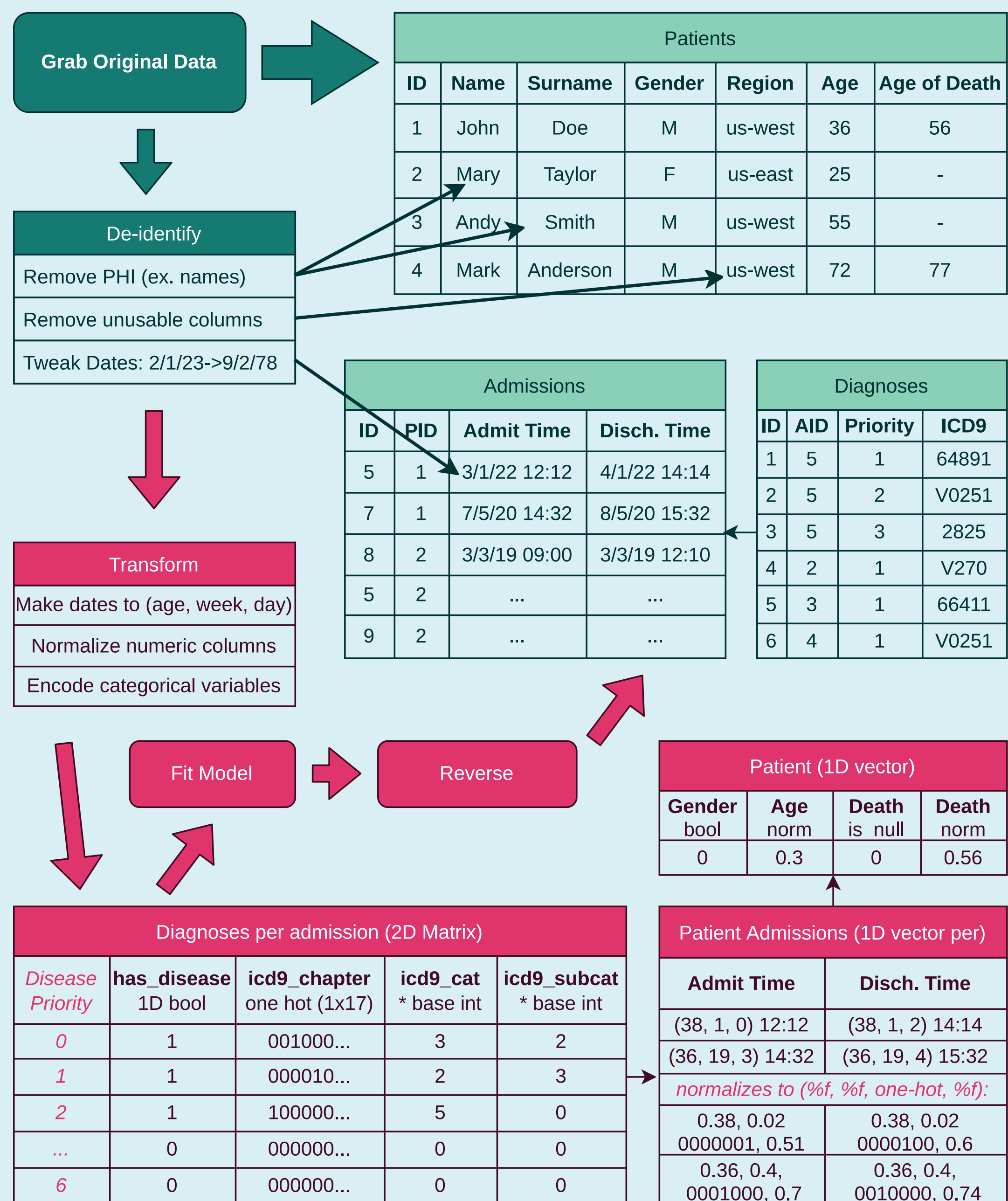
First, when modeling the data, great care is taken to create a **compressed representation** (a dense **bayesian network** or a small **neural network** with transfer learning). This limits the information extracted from the data, lowering the **noise** added. Then, **Differential Privacy** is applied to the model, which mathematically bounds the privacy and accuracy output.

The end result is **benchmarked** by measuring the **deviation of queries**, the behavior of **models**, and a suite of **marginal metrics** (statistics based).

## Future goals

Test viability of **Differential Privacy** in **Data Synthesis**, create multiple models for **synthesizing multimodal data**, and create a suite of **viable metrics** for that data.

## Privacy Aware Data Synthesis



## Testing of Synthetic Data

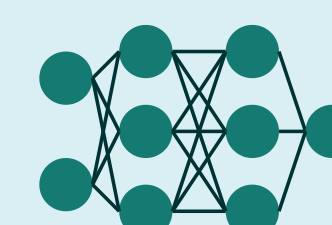
CS	$\chi^2$	p
lang	0.01	0.89
ethn	0.03	0.7


KS	KS	p
age	0.02	0.78
death	0.05	0.68

KL	diag	meds
age	0.8	0.78
death	0.7	0.5

**AVG err:** `select avg(age) from patients`  
`select count(*) from diagnoses group by icd9 limit 20`

## Train classifiers, Test accuracy:



$\uparrow \text{accuracy}_{\text{real data}} = \uparrow \text{quality}$    
 $\text{accuracy}_{\text{train data}} > \text{accuracy}_{\text{test data}} = \downarrow \text{privacy}$