

# Health Data Anonymization in Theory and Practice

Prof. Dr. Fabian Prasser  
Medical Informatics Group  
Berlin Institute of Health @  
Charité – Universitätsmedizin Berlin

# Motivation: Data Sharing and Re-Use

- **Data-driven approaches in medical research**
  - Precision medicine: high case numbers, detailed characterizations
  - Real-world evidence: secondary use, e.g. of routine clinical data for research
  - Collaborative research, e.g. data sharing across institutional boundaries
- **Initiatives to improve the transparency, reproducibility and reusability of research results and research data**
  - NIH Statement on Sharing Research Data, Notice NOT-OD-03-032; 2003.
  - NIH Genomic Data Sharing Policy, Notice NOT-OD-14-124; 2014.
  - EMA Policy 0070 on Publication of Clinical Data for Medicinal Products for Human Use; 2014.
  - Increased citation rates

# Balancing Interests When Sharing Data Can Be Difficult



Transparency



Innovation



Common  
good



Ethics



Privacy

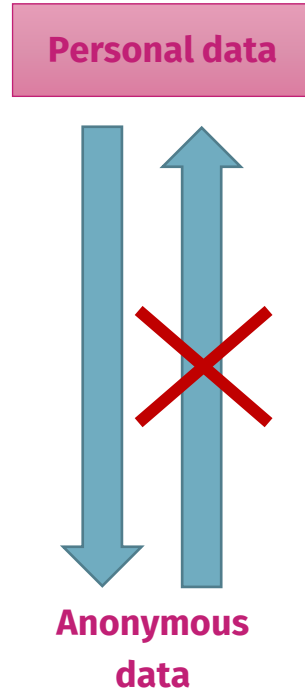


Regulations



Intellectual  
property

# EU General Data Protection Regulation (GDPR)



GDPR, Recital 26:

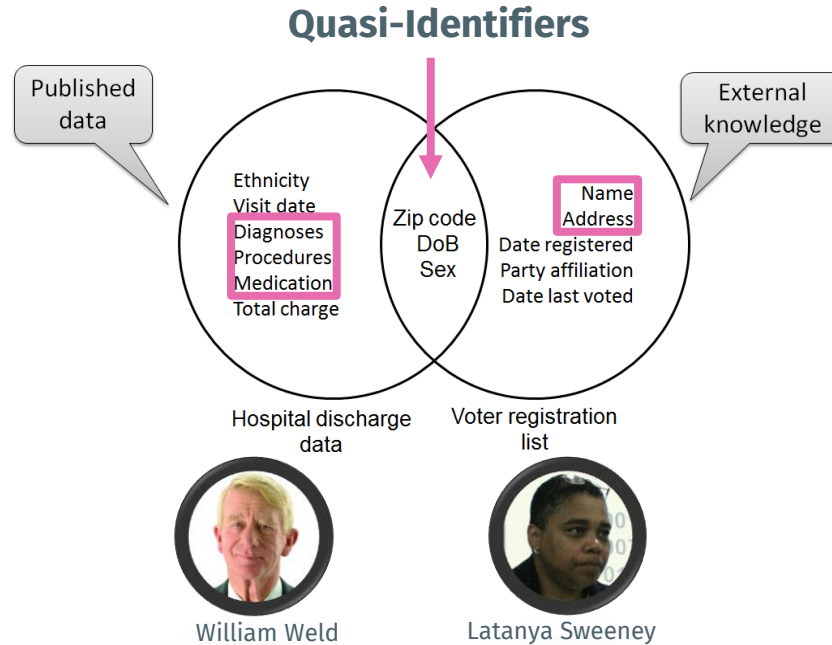
„The principles of data protection should **apply to any information concerning an identified or identifiable natural person** [...]“

„[...] To determine whether a natural person is identifiable, **account should be taken of all the means reasonably likely to be used**, [...] to identify the natural person directly or indirectly [...]“

"[In doing so] all **objective factors**, such as the costs of and the **amount of time required** for identification, taking into consideration the **available technology at the time of the processing and technological developments** [...]"

Source: Regulation (EU) 2016/679 of the European parliament and the council of 27 April 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC (General Data Protection Regulation)

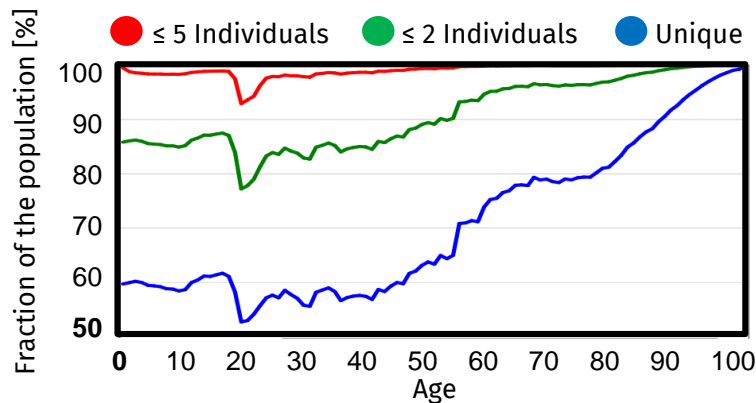
# The Case of William Weld (1997): Linkage Attacks



Sources: Golle P. Revisiting the uniqueness of simple demographics in the US population. 5th ACM Workshop on Privacy in the Electronic Society, 2006, Sweeney L. Simple Demographics Often Identify People Uniquely. Carnegie Mellon University, Data Privacy Working Paper 3. Pittsburgh 2000, Image by Gary Johnson from Taos, NM - BillWeld5x7 (2), CC BY 2.0, <https://commons.wikimedia.org/w/index.php?curid=49683363>

# Uniqueness of Simple Demographics

87 % of the US population can be uniquely identified by the combination of DoB, Sex, ZIP Code



Sources: Golle P. Revisiting the uniqueness of simple demographics in the US population. 5th ACM Workshop on Privacy in the Electronic Society, 2006, Sweeney L. Simple Demographics Often Identify People Uniquely. Carnegie Mellon University, Data Privacy Working Paper 3. Pittsburgh 2000.

# Re-Identification Revisited (2019)

## Medical Data De-Identification Is Under Attack



David Talby Forbes Councils Member  
Forbes Technology Council COUNCIL POST | Paid Program  
Innovation

POST WRITTEN BY

David Talby

PhD, MBA, CTO at [Pacirc AI](#). Making AI, big data and data science solve real-world problems in healthcare, life science and related fields.

Forbes - Forbes Technology Council, 27.08.2019

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.

The New York Times, 23.07.2019

## “Anonymous” Data Won’t Protect Your Identity

A new study demonstrates it is surprisingly easy to ID an individual within a supposedly incognito data set

Scientific American, 23.07.2019



ARTICLE

<https://doi.org/10.1038/s41467-019-10933-3> OPEN

Estimating the success of re-identifications in incomplete datasets using generative models

Luc Rocher<sup>1,2,3</sup>, Julien M. Hendrickx<sup>1</sup> & Yves-Alexandre de Montjoye<sup>2,3</sup>

Nature Communications, 23.07.2019

“[...] we find that 99.98% of Americans would be correctly re-identified in any dataset using 15 demographic attributes.”

# Any Characteristic Can Make You Unique!

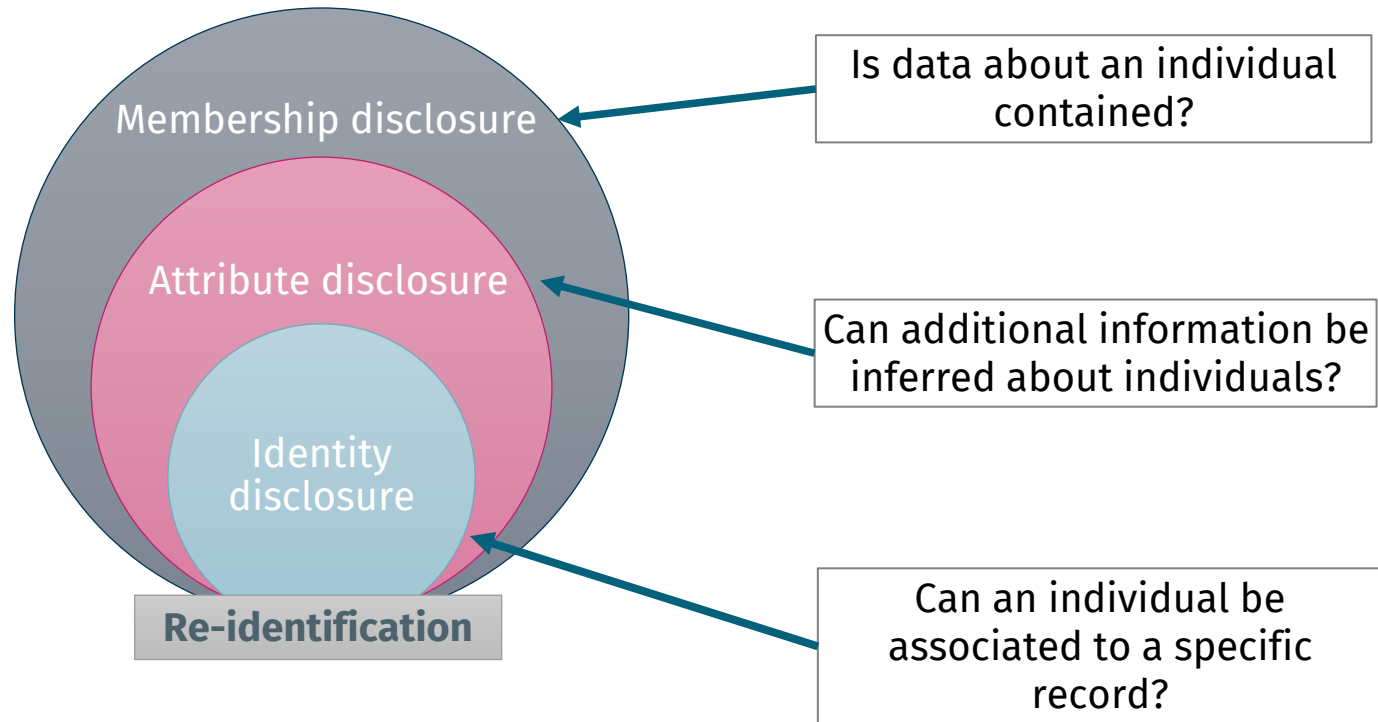
- **Demographic data** (Sweeney 1997; Golle 2006; El Emam 2008)
- **Diagnosis codes** (Loukides et al. 2010)
- **DNA (SNPs)** (Lin, Owen, & Altman 2004; Homer et al. 2008, Wang et al. 2009)
- **Pedigree structure** (Malin 2006)
- **Location visits** (Malin & Sweeney 2004, Golle & Partridge 2009)
- **Movie reviews** (Narayanan & Shmatikov 2008)
- **Search queries** (Barbaro & Zeller 2006)
- **Social network structure** (Backstrom et al. 2007, Narayanan & Shmatikov 2009)

But: Unique  $\neq$  Identifiable!

Source: Bradley Malin. Challenges and Solutions for Data Privacy in Translational Research. 2011



# (Re-)Identification: Types of Disclosure



# (Re-)Identification: Example

Direct Id.		Quasi-Identifiers			Ins.	Sensitive Information
Name	Tel.	Age	Sex	ZIP	BMI	Diagnosis
*	*	[50,60)	M	92***	28.0	I47.1 Supraventricular tachycardia
*	*	[90,100)	M	94***	24.9	I50.1 Left ventricular failure
*	*	[80,90)	F	92***	26.7	I50.0 Congestive heart failure
*	*	[60,70)	F	94***	31.7	I47.1 Supraventricular tachycardia
*	*	[70,80)	F	92***	18.3	I47.0 Re-entry ventricular arrhythmia
*	*	[80,90)	F	92***	24.0	I50.0 Congestive heart failure
*	*	[80,90)	F	92***	28.1	I50.0 Congestive heart failure
*	*	[70,80)	M	94***	31.0	I47.0 Re-entry ventricular arrhythmia
*	*	[80,90)	M	94***	34.9	I50.0 Congestive heart failure
*	*	[60,70)	M	93***	32.3	I50.1 Left ventricular failure

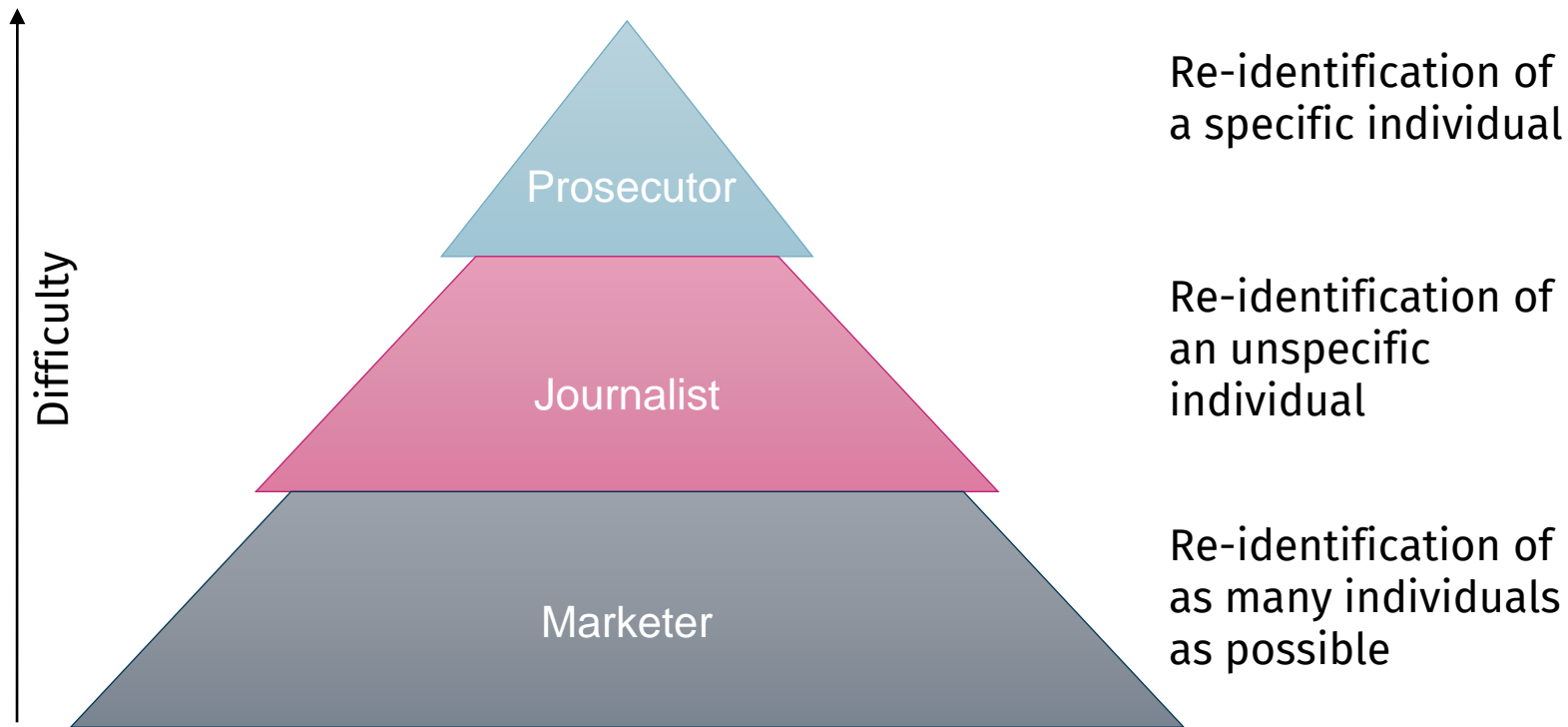
1

2

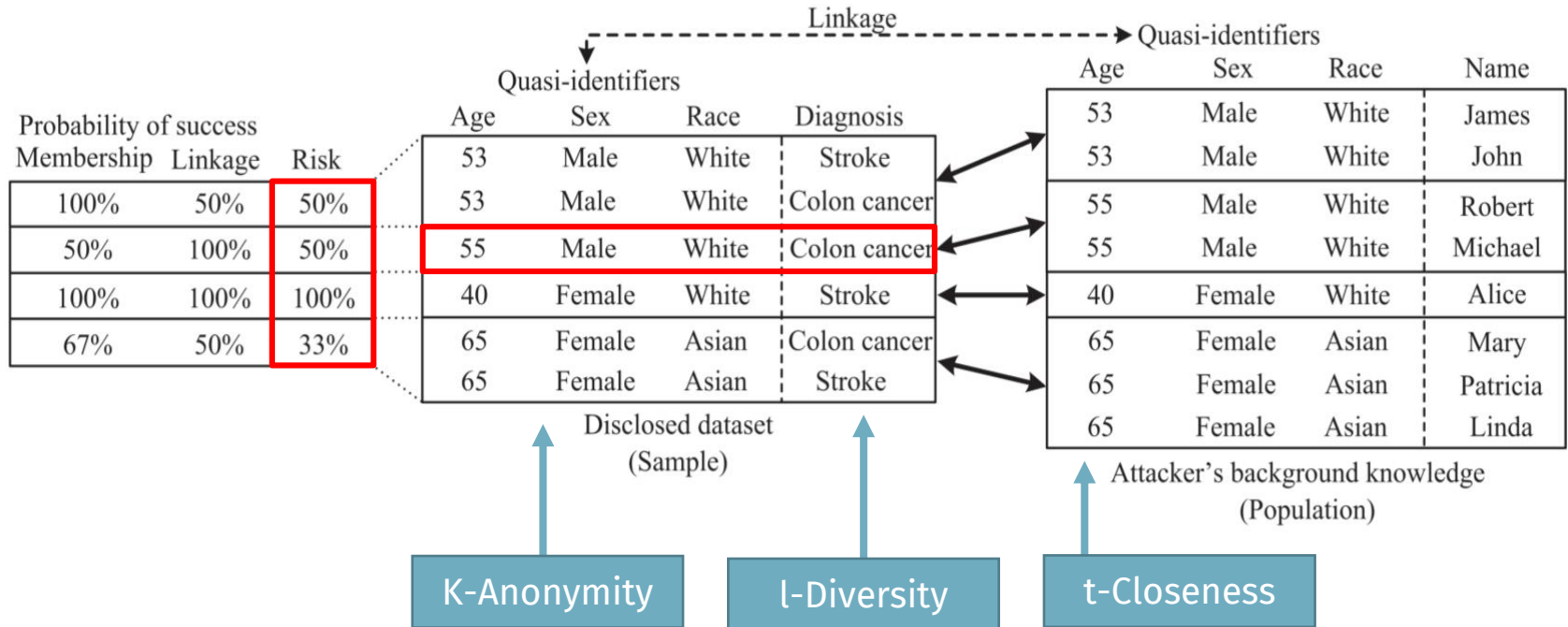
3

- (1) Membership Disclosure
- (2) Attribute Disclosure
- (3) Identity Disclosure

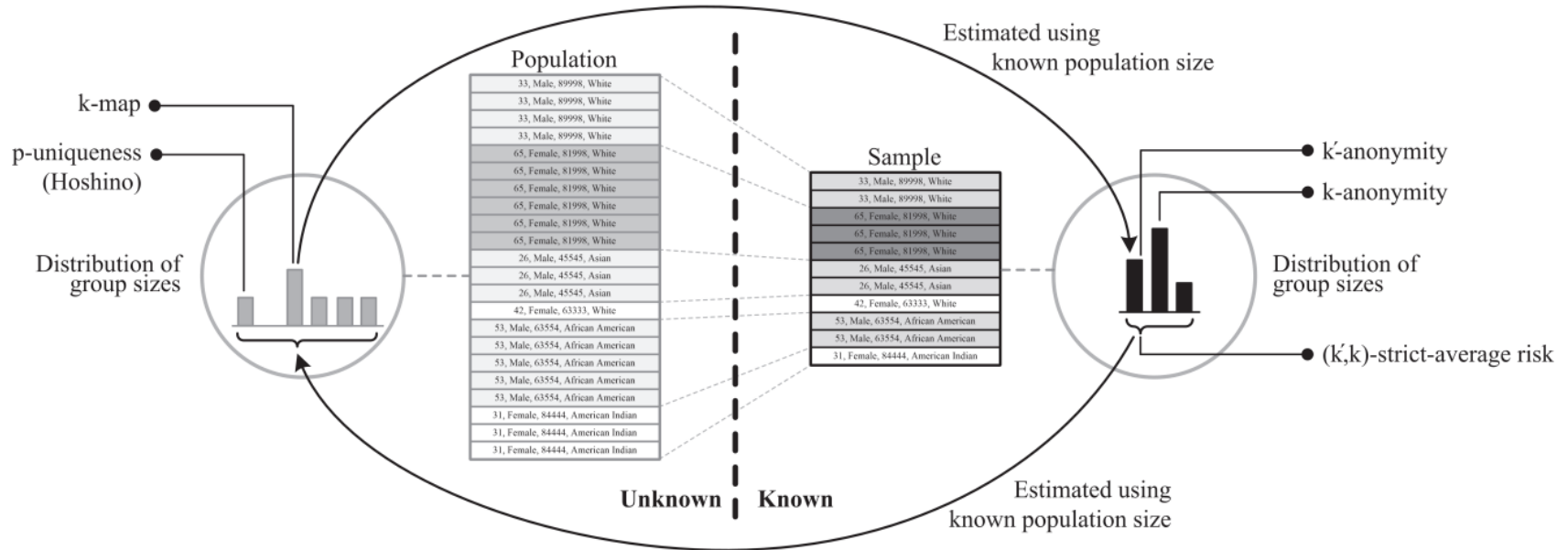
# (Re-)Identification: Attacker Goals



# Examples of Risk Models



# Examples of Risk Models (2)



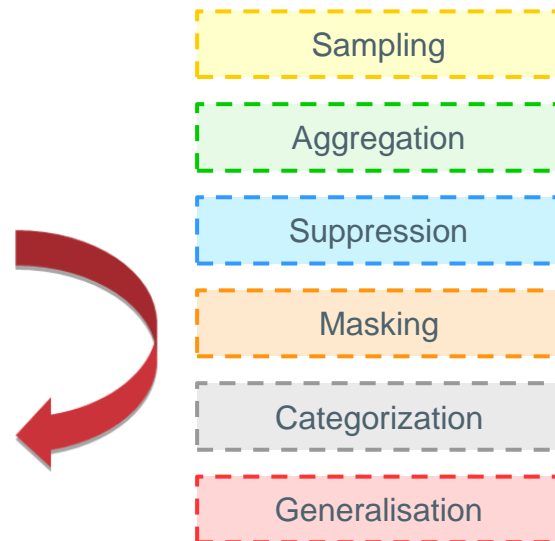
# Anonymization: A Basic Example

Processing of personal (input) data in such a way that anonymous (output) data is produced.

Example:

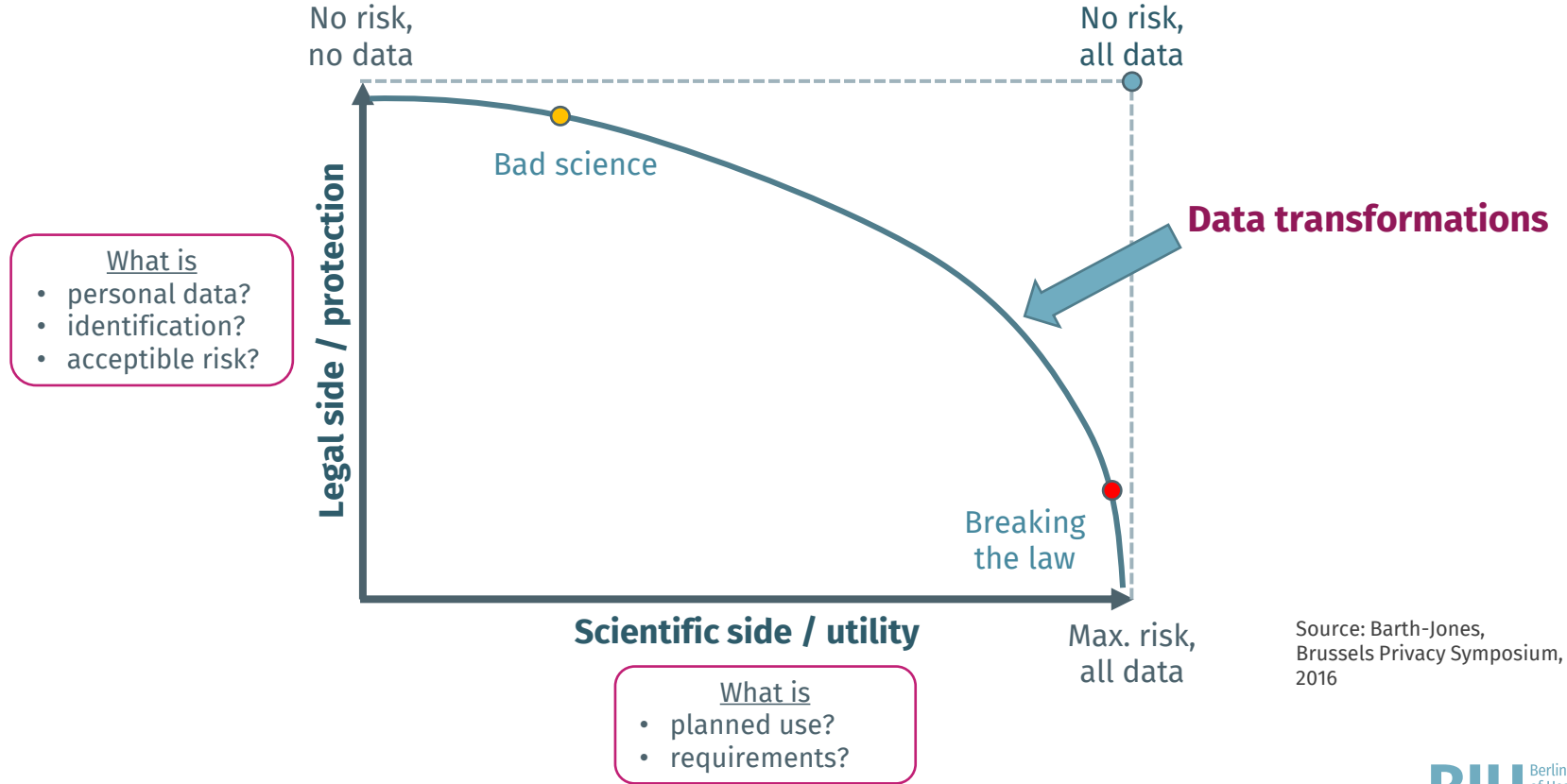
Alter	Geschlecht	PLZ	Gewicht	Diagnose
55	Männlich	81539	71	C25.0 Bösartige Neubildung des Pankreas - Pankreaskopf
76	Männlich	81675	80	C25.0 Bösartige Neubildung des Pankreas - Pankreaskopf
66	Männlich	81929	85	C25.0 Bösartige Neubildung des Pankreas - Pankreaskopf
81	Männlich	80802	79	C25.1 Bösartige Neubildung des Pankreas - Pankreaskörper
74	Männlich	81249	88	C25.2 Bösartige Neubildung des Pankreas - Pankreasschwanz
71	Weiblich	80335	69	C18.2 - Bösartige Neubildung des Kolons - Colon ascendens
64	Weiblich	80339	71	C18.4 - Bösartige Neubildung des Kolons - Colon transversum
69	Männlich	80637	75	C18.7 - Bösartige Neubildung des Kolons - Colon sigmoideum
55	Weiblich	80638	77	C18.7 - Bösartige Neubildung des Kolons - Colon sigmoideum
61	Männlich	81667	67	C18.7 - Bösartige Neubildung des Kolons - Colon sigmoideum

Alter	Geschlecht	PLZ	Gewicht	Diagnose
72,0	Männlich	81***	[80, 90[	C25.- Bösartige Neubildung des Pankreas
72,0	Männlich	81***	[80, 90[	C25.- Bösartige Neubildung des Pankreas
72,0	Männlich	81***	[80, 90[	C25.- Bösartige Neubildung des Pankreas
62,7	---	80***	[70, 80[	C18.- Bösartige Neubildung des Kolons
62,7	---	80***	[70, 80[	C18.- Bösartige Neubildung des Kolons
62,7	---	80***	[70, 80[	C18.- Bösartige Neubildung des Kolons



**k-Anonymity and ( $\epsilon$ ,  $\delta$ )-Differential Privacy**

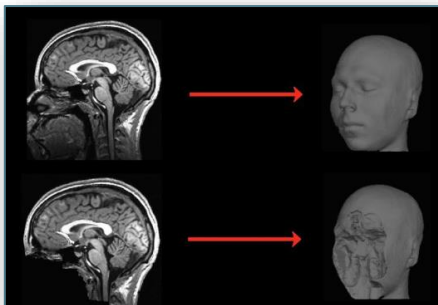
# Trade-Offs in Anonymization



Source: Barth-Jones, Brussels Privacy Symposium, 2016

# Risk-Based Anonymization

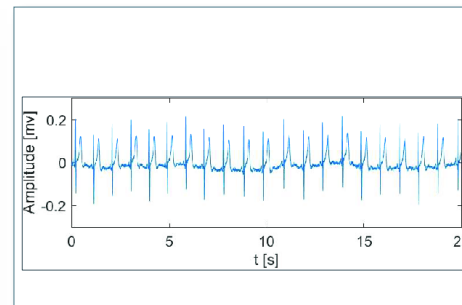
Context: Purpose, recipient, types of data etc.



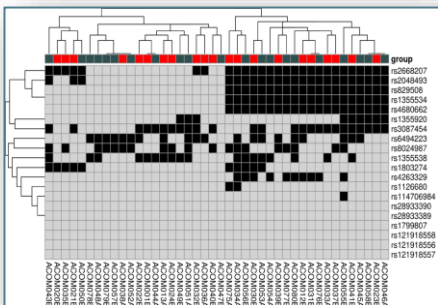
Source: [https://surfer.nmr.mgh.harvard.edu/fswiki/mri\\_deface](https://surfer.nmr.mgh.harvard.edu/fswiki/mri_deface)

AUTOPSY REPORT - Final Anatomic Diagnosis  
 Dx: Sickle cell anemia with multiple red blood cell transfusions  
 Cause of death per autopsy report (ICD-9-CM): Cirrhosis related to Hepatitis B  
 Mr. **Norman Hesse** is a 50 year old male, originally from **Dr. Leona**, who was diagnosed with sickle cell anemia at age 8. From the age of 8 to 18, he had several health complications and underwent a liver transplant at the **Casevik Hospital Center** in **Michigan** 2014. He has been in good health and continued with normal daily activities until **June 2016**, when he was brought to the **Steppenwolf Clinic** and admitted to the **ICU**. At that time, he was diagnosed with end-stage renal disease. He responded well to hemodialysis for about a year per his **Dr. Bernice Morari**. A few months later he began to experience chronic pain in his left hip and was referred to Dr. **Joshua** at the **Steppenwolf Well Pain Management Center**. On **October 14th, 2017**, he was re-admitted to the **Steppenwolf Clinic** and quickly transferred to the **ICU**. Due to his declining health, the patient's **Dr.** met with an **ethics consultant** and decided to withdraw medical services and provide comfort measures only. The patient expired on **October 18th, 2017**. A limited autopsy was performed on the **stomach** at **11:00am**.

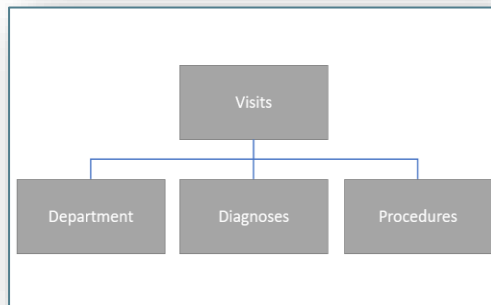
Source: <https://scrubber.nlm.nih.gov/>



Source: <https://doi.org/10.1109/MeMeA.2018.8438751>



Source: <https://doi.org/10.2147/CCID.S176842>



Source: [https://www.g-drg.de/Datenlieferung\\_gem\\_21\\_KHEntgG](https://www.g-drg.de/Datenlieferung_gem_21_KHEntgG)

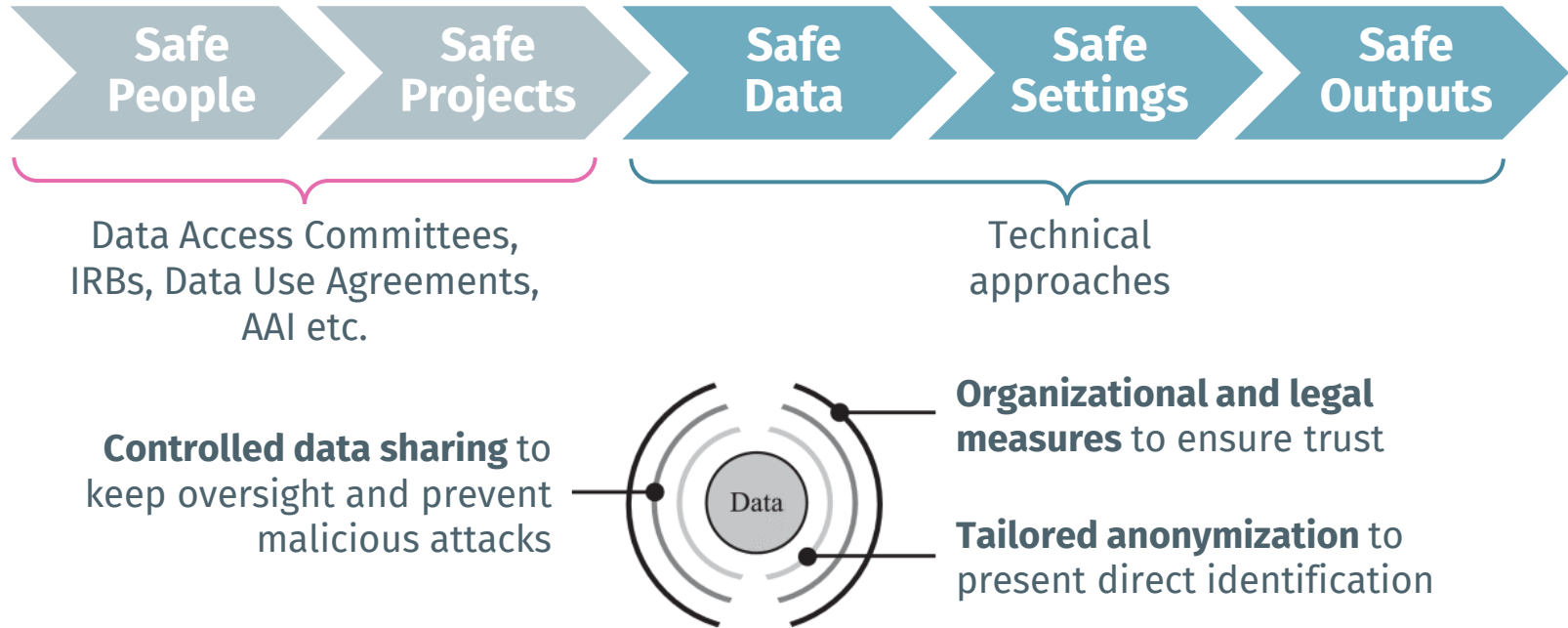
Onset of exposure	Yes	No	Total
20+ years***	339	53	392
0-19 years***	203	522	725
Total	542	575	1,117

Source: <https://doi.org/10.1080/10937404.2012.678766>



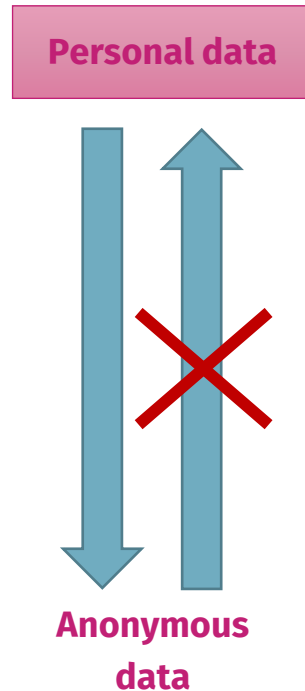
# Controlling the Context

## The Five Safes Framework



Source: Desai, Ritchie, Welpton. (2016) Five Safes: designing data access for research.

# Risk-Based Anonymization and the GDPR



GDPR, Recital 26:

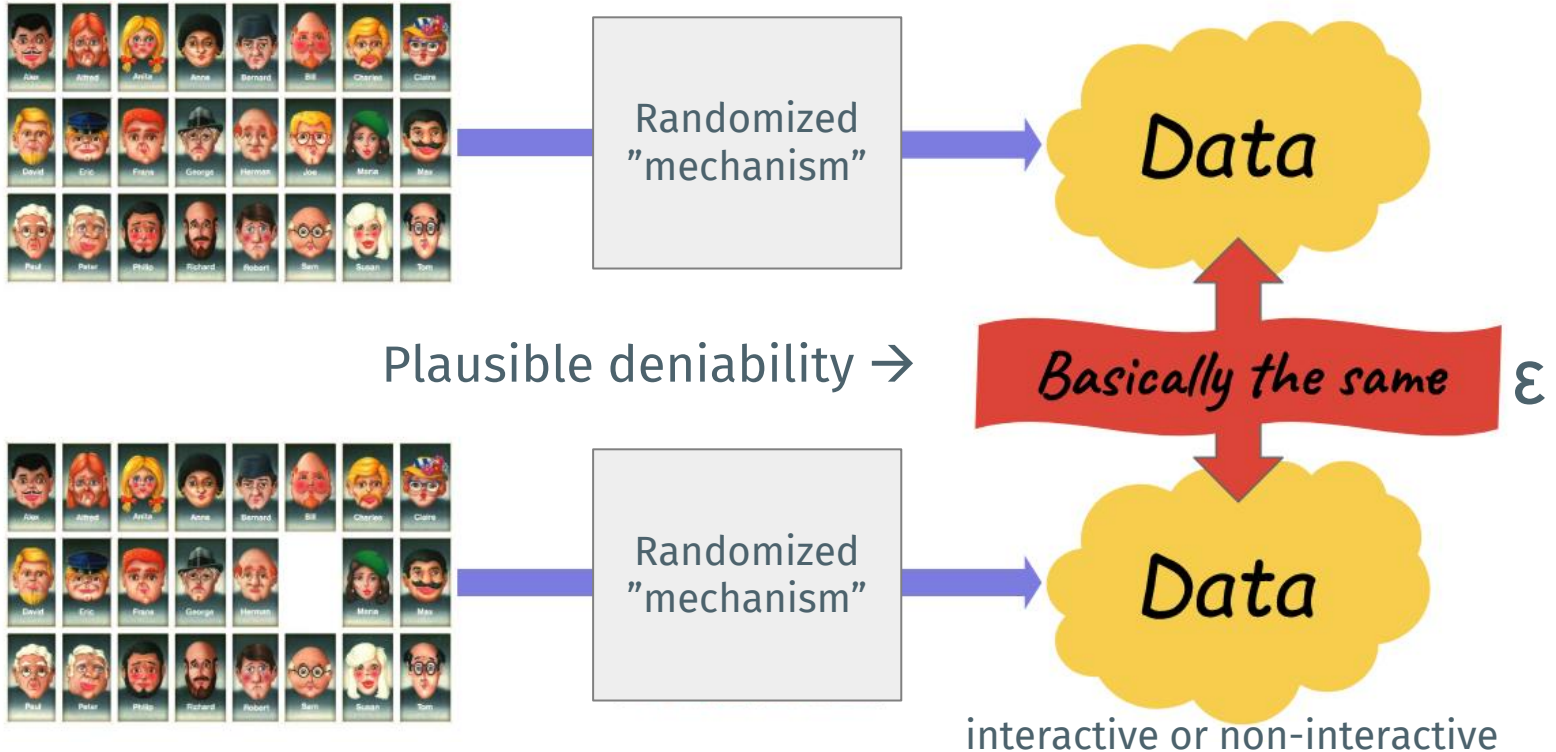
„The principles of data protection should **apply to any information concerning an identified or identifiable natural person** [...]“

„[...] To determine whether a natural person is identifiable, **account should be taken of all the means reasonably likely to be used**, [...] to identify the natural person directly or indirectly [...]“

"[In doing so] all **objective factors**, such as the **costs of and the amount of time required** for identification, taking into consideration the **available technology at the time of the processing and technological developments** [...]"

Source: Regulation (EU) 2016/679 of the European parliament and the council of 27 April 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC (General Data Protection Regulation)

# A New Perspective: Differential Privacy



Source: <https://desfontain.es/privacy/differential-privacy-in-more-detail.html>

# Overview of Available Tools

Tool	Institution	Country	Language(s)	Release	Last Update	License
<b>μ-Argus</b>	Centraal Bureau voor de Statistiek	Netherlands	C++, Java	1998	2021	EUPL
<b>sdcmicro</b>	Statistics Austria	Austria	R	2007	2021	GPL 2
<b>Open Anonymizer</b>	University of Vienna	Austria	Java	2008	2009	Unknown
<b>CAT</b>	Cornell University	USA	C++	2009	2014	Unknown
<b>Tiamat</b>	Purdue University	USA	Java	2009	Unknown	Unknown
<b>UTD</b>	The University of Dallas	USA	Java	2010	2012	GPL 2
<b>Anon</b>	University of Klagenfurth	Austria	Java	2012	Unknown	Unknown
<b>ARX</b>	BIH@Charité	Germany	Java	2012	2022	Apache 2
<b>SECRET</b>	University of Peloponnes	Greece	C++, Qt	2013	Unknown	Unknown
<b>Probabilistic Anonymization</b>	University of Cyprus, Cyprus and Newcastle University, UK	Greece/UK	R	2018	2018	Unknown
<b>μ-Ant</b>	Center for Cybersecurity Research of Catalonia	Spain	Java	2019	2019	MIT
<b>Amnesia</b>	University of Thessaly	Greece	Java, JavaScript	2019	2022	BSD 3-Clause
<b>PrioPrivacy</b>	Research Studio Data Science	Austria	Java	2019	2021	Unknown

- Time focus around 2010 and 2020
- Most tools come from European institutions
- Most common programming languages are Java, C++, R
- Half of the tools identified are only publicly accessible to a limited extent
- Only a few tools are under permanent development
- Further research shows that only three of the tools (μ-Argus, ARX, sdcmicro) are used in real application scenarios.

# Privacy Models Supported by Tools

Privacy Model	Model Type	$\mu$ -Argus	sdcmicro	Open Anonymizer	CAT	TIAMAT	UTD	Anon	ARX	SECRET	$\mu$ -Ant	Amnesia	PrioPrivacy	Probabilistic Anonymization
<b>k<sup>m</sup>-Anonymity</b>	Syntactic/statistical									X		X		
<b>k-Anonymity</b>	Syntactic/statistical	X*	X	X		X	X	X	X	X	X	X	X	
<b>ℓ-Diversity</b>	Syntactic/statistical				X		X	X	X					
<b>t-Closeness</b>	Syntactic/statistical				X		X		X		X			

\*Only for combinations of one, two or three key variables

- Restriction to models that (1) can be implemented automatically by the tools and (2) are supported by at least two tools.
- Amnesia and SECRET are the only tools that support a model for set-valued attributes.
- Only four tools offer models for protection against attribute disclosure

# sdcmicro

## Developed at Statistics Austria / Technische Universität Wien

- Under development since 2008

## Primary application domain

- Disclosure control in official statistics

## Methods

- Focus on the a posteriori methodology
- Highly flexible because of R integration

## Interfaces

- Primarily a package for the R statistics software
- Cross-platform graphical user interface

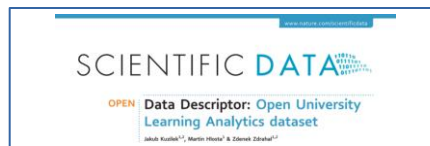
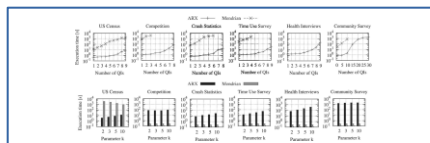
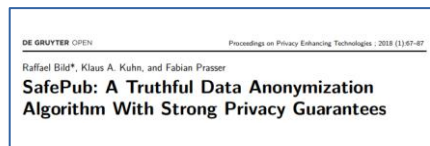
## Scalability

- Highly scalable when used as a programming library
- Scalability issues when using the graphical interface



# ARX: Features and Applications

- **Comprehensive feature set:** „traditional“ approaches, Differential Privacy, game-theoretic methods, privacy-preserving machine learning.
- **Quite scalable:** Significantly outperforms related tools, used to anonymise datasets with billions of records.
- **Graphical tool:** Used in education and training by commercial and public institutions in several countries.
- **Wide range of applications:** Creation of open datasets and used to build anonymisation pipelines in several domains, e.g. by telecom providers, health insurances.
- **Industry friendly:** Integrated into several commercial products, core algorithms adopted by SAP HANA.
- **Open source:** More than 50.000 downloads.



# ARX: Graphical Frontend

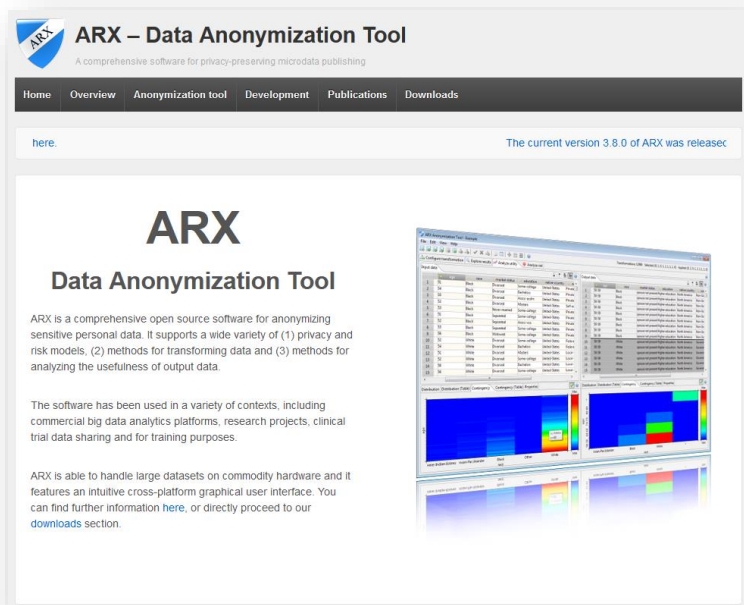
The image displays the ARX graphical frontend, a tool for data anonymization. It features several interconnected windows:

- Configuration:** A window for setting transformation rules, such as "Quasi-identifying" and "Generalization". It includes a table for defining levels of generalization (Level-0 to Level-4) and options for privacy criteria and suppression.
- Exploration:** A window showing a large grid of data points, likely representing the results of a transformation or a search space exploration.
- Quality analysis:** A window displaying summary statistics and contingency tables, used to evaluate the quality of the anonymized data.
- Risk analysis:** A window showing a histogram of prosecutor re-identification risk, with a risk threshold line and a distribution of records affected by different risk levels.

A central circular diagram with four quadrants (Configuration, Exploration, Risk analysis, Quality analysis) and arrows indicates a continuous cycle between these four main components of the tool.

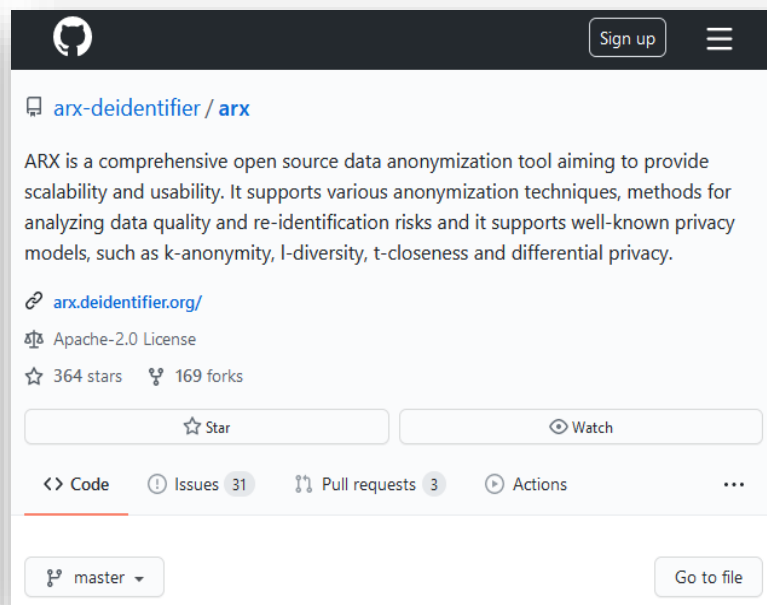


# ARX: Open Source Project



The screenshot shows the ARX website homepage. At the top left is the ARX logo and the title "ARX – Data Anonymization Tool" with the subtitle "A comprehensive software for privacy-preserving microdata publishing". A navigation bar includes links for Home, Overview, Anonymization tool, Development, Publications, and Downloads. A message states "The current version 3.8.0 of ARX was released here." The main content area features the title "ARX Data Anonymization Tool" and a description: "ARX is a comprehensive open source software for anonymizing sensitive personal data. It supports a wide variety of (1) privacy and risk models, (2) methods for transforming data and (3) methods for analyzing the usefulness of output data." Below this, it mentions the software's use in commercial big data analytics, research, and clinical trial data sharing. A screenshot of the ARX software interface is shown, displaying a data table and a heatmap. At the bottom, it states that ARX can handle large datasets on commodity hardware and features an intuitive cross-platform graphical user interface, with links to find further information and the downloads section.

Comprehensive website

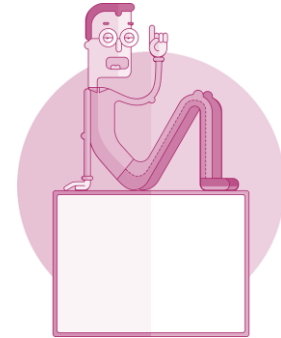


The screenshot shows the ARX GitHub repository page. At the top right, there is a "Sign up" button and a menu icon. The repository name is "arx-deidentifier / arx". The description states: "ARX is a comprehensive open source data anonymization tool aiming to provide scalability and usability. It supports various anonymization techniques, methods for analyzing data quality and re-identification risks and it supports well-known privacy models, such as k-anonymity, l-diversity, t-closeness and differential privacy." Below the description, there is a link to "arx.deidentifier.org/", the license "Apache-2.0 License", and statistics showing "364 stars" and "169 forks". There are buttons for "Star" and "Watch". At the bottom, there are links for "Code", "Issues 31", "Pull requests 3", and "Actions". A dropdown menu shows "master" and a "Go to file" button is present.

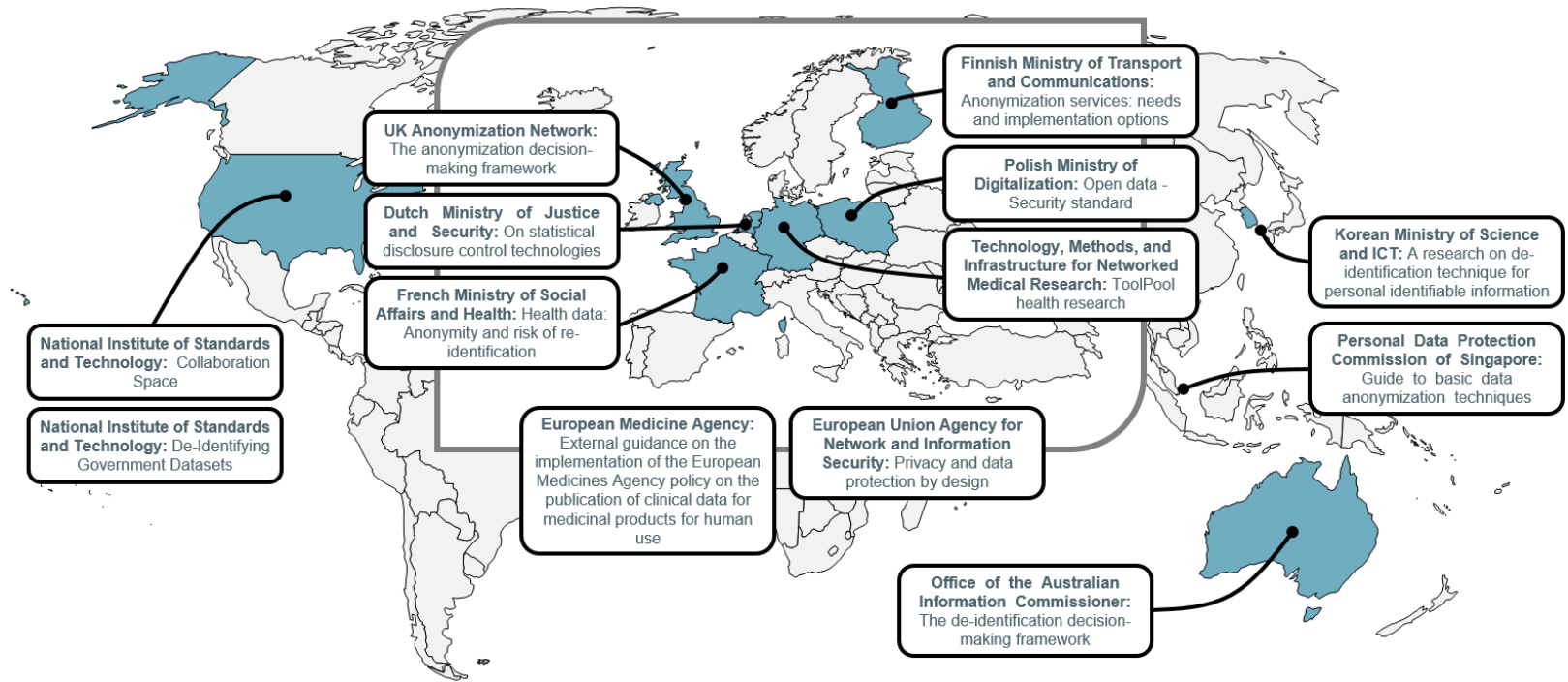
Community contributions

# Examples of Guidelines Mentioning ARX (1)

- European Medicines Agency. EMA/90915/2016 – external guidance on the implementation of the European medicines agency policy on the publication of clinical data for medicinal products for human use; 2018.
- European Union Agency for Network and Information Security. Privacy and data protection by design; 2015.
- UKAN. The anonymisation decision-making framework; 2016.
- Office of the Australian Information Commissioner. The de-identification decision-making framework; 2017.
- French Ministry of Solidarity and Health. Health data: anonymity and risk of re-identification; 2015.
- Finnish Ministry of Transport and Communications. Anonymization services – requirements and implementation options; 2017.
- Personal Data Protection Commission of Singapore. Guide to basic data anonymisation techniques; 2018.
- Polish Ministry of Digitalization. Open data - Security standard; 2018.
- Dutch Ministry of Justice and Security. On statistical disclosure control technologies; 2018.
- Korean Ministry of Science and ICT. A research on de-identification technique for personal identifiable information; 2016.



# Examples of Guidelines Mentioning ARX (2)



World Map provided by simplemaps.com

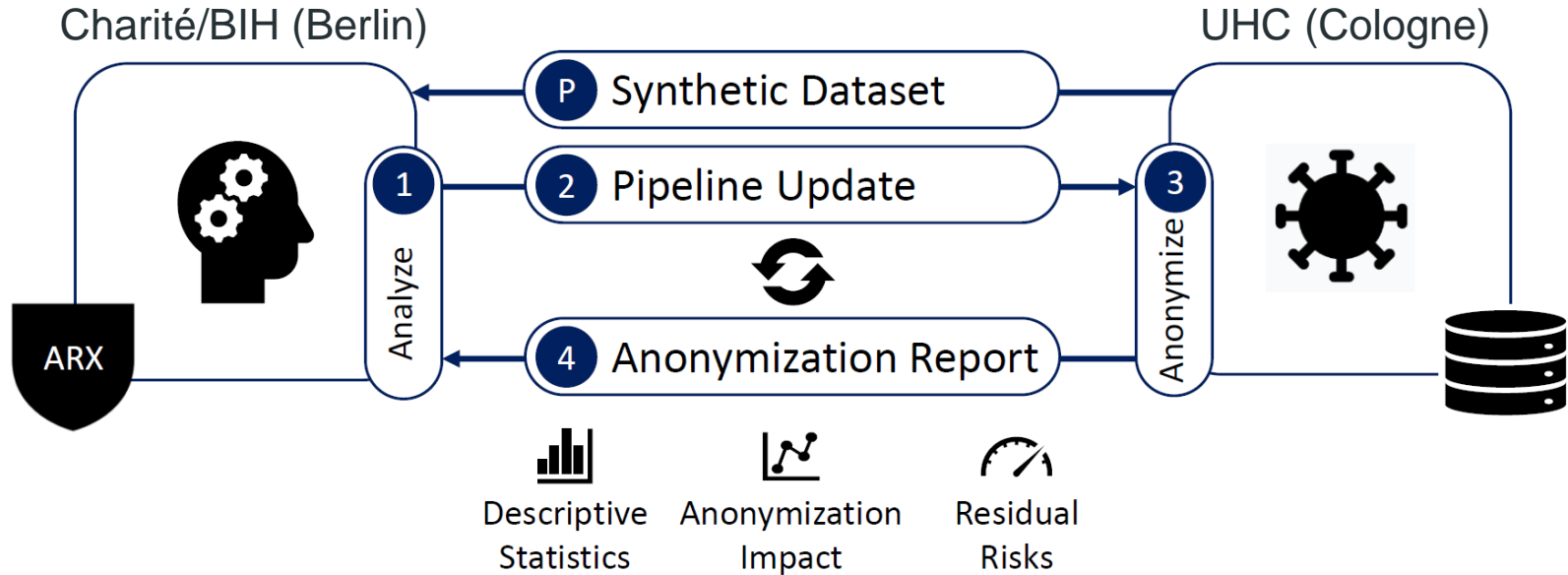
# Example: Anonymisation Pipelines for the LEOSS registry

- **LEOSS: A European registry capturing the clinical course of SARS-CoV-2 infected patients (<https://leoss.net>) established at University of Cologne**
  - No informed consent necessary (anonymous reports).
  - Retrospective documentation after discharge / death.
  - All hospitalized patients including children eligible.
  - Immediate start after verification.
- **Open Science approach**
  - Registry hosted in a secure environment in Cologne.
  - Anonymous data is shared with researchers and the public.
  - Additional anonymisation procedures have been implemented for this purpose.

# LEOSS: Overview

- **Two types of datasets**
  - Public Use File with 16 variables available without restrictions.
  - Scientific Use Files with  $\leq 605$  variables available under data use contracts.
- **Two types of pipelines, built with ARX**
  - Two stages for the Public Use File
  - Ten stages for the Scientific Use File
- **Both pipelines were developed without access to primary data in close cooperation with the LEOSS Core Team in Cologne.**

# LEOSS: Development Process



→ Seven iterations over several weeks

# LEOSS: Approach for the Public Use File (1)

## (1) Qualitative risk assessment

- Compared data to “risky” variables mentioned in laws and guidelines.
  - Low risk already according to this initial assessment.
- Additionally, assessed the risk of identification associated with individual variables following a methodology proposed by Malin et al.\*
  - Replicability, availability, distinguishability categorized into low, medium or high.
  - Variables above threshold considered potentially identifying.

## (2) Quantitative risk assessment

- Followed recommendations from the Opinion on Anonymisation Methods by the Article 29 Data Protection Working Party (today: European Data Protection Board):
  - Singling out: the possibility to isolate some or all records which identify an individual in the dataset.
  - Linkability: the ability to link, at least, two records concerning the same data subject or a group of data subjects.
  - Inference: the possibility to deduce, with significant probability, the value of an attribute from the values of a set of other attributes.



\* Malin, B., Loukides, G., Benitez, K. & Clayton, E. W. Identifiability in biobanks: models, measures, and mitigation strategies. Hum. Genet. 130, 383–392 (2011).

# LEOSS: Approach for the Public Use File (2)

## (3) Formal anonymization process

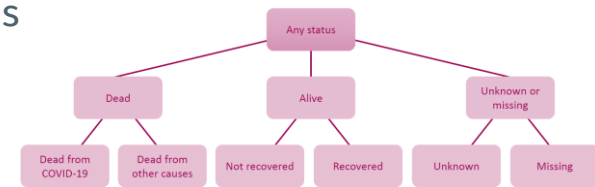
- Generalization and record suppression to mitigate risks highlighted by the Opinion.
- Prevented singling out and linkability by reducing the uniqueness of all possible combinations of potentially identifying variables (k-anonymity).
- Prevented inference by ensuring that the distribution of medical data within groups of indistinguishable records is not too different from the distribution in the overall dataset (t-closeness).
- Static generalization scheme and withholding of records to ensure that protection holds also when data is updated repeatedly.

## (4) Extensive documentation

- Entire development process and underlying considerations are documented in detail. Pipeline released as OSS.

## (5) Continuous monitoring

- Repeated evaluation of data utility.





# LEOSS: Result

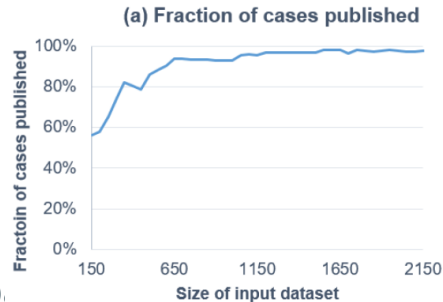
Variable	Description
Age at diagnosis	Age of patient at time of diagnosis
Gender	Sex of patient
Month first diagnosis	Month of first confirmed diagnosis of COVID-19
Year first diagnosis	Year of first confirmed diagnosis of COVID-19
Uncomplicated phase	Indicates whether the patient has been through the uncomplicated phase of COVID-19
Complicated phase	Indicates whether the patient has been through the complicated phase of COVID-19
Critical phase	Indicates whether the patient has been through the critical phase of COVID-19
Recovery phase	Indicates whether the patient has been through the recovery phase of COVID-19
Vasopressors in complicated phase	Indicates whether vasopressors were used in the complicated phase
Vasopressors in critical phase	Indicates whether vasopressors were used in the critical phase
Invasive ventilation in critical phase	Indicates whether invasive ventilation was used in the critical phase
Superinfection in uncomplicated phase	Type of (if any) superinfection in uncomplicated phase
Superinfection in complicated phase	Type of (if any) superinfection in complicated phase
Superinfection in critical phase	Type of (if any) superinfection in critical phase
Symptoms in recovery phase	Symptoms (if any) in recovery phase
Last known patient status	Last known status

# LEOSS: Evaluation (1)

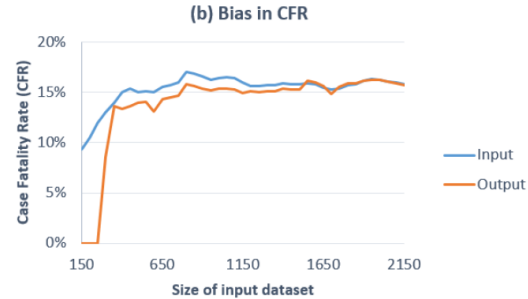
## Pipeline based on the principle of “hiding in the crowd”

- Anonymity is achieved by making sure that each record does not differ significantly from a larger group of records.
- Counter-intuitive property: the greater the number of individuals included in the registry, the less information has to be removed to achieve the required degree of protection.

## Example: records released and case fatality rate

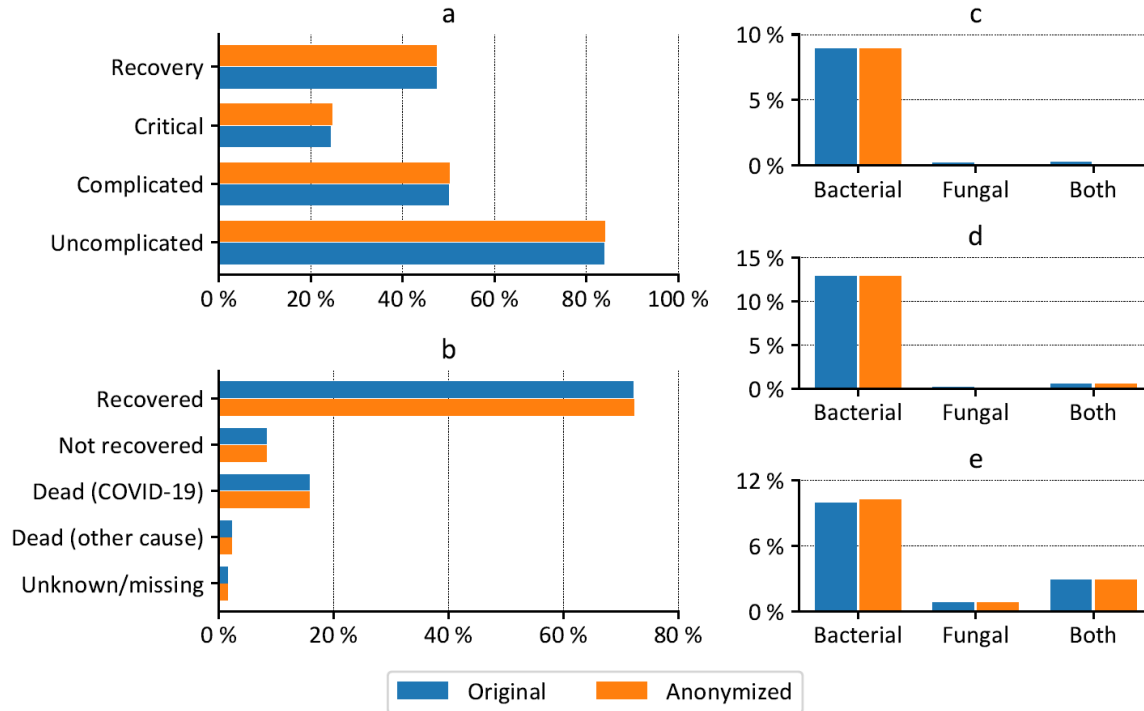


→ Negligible imp.



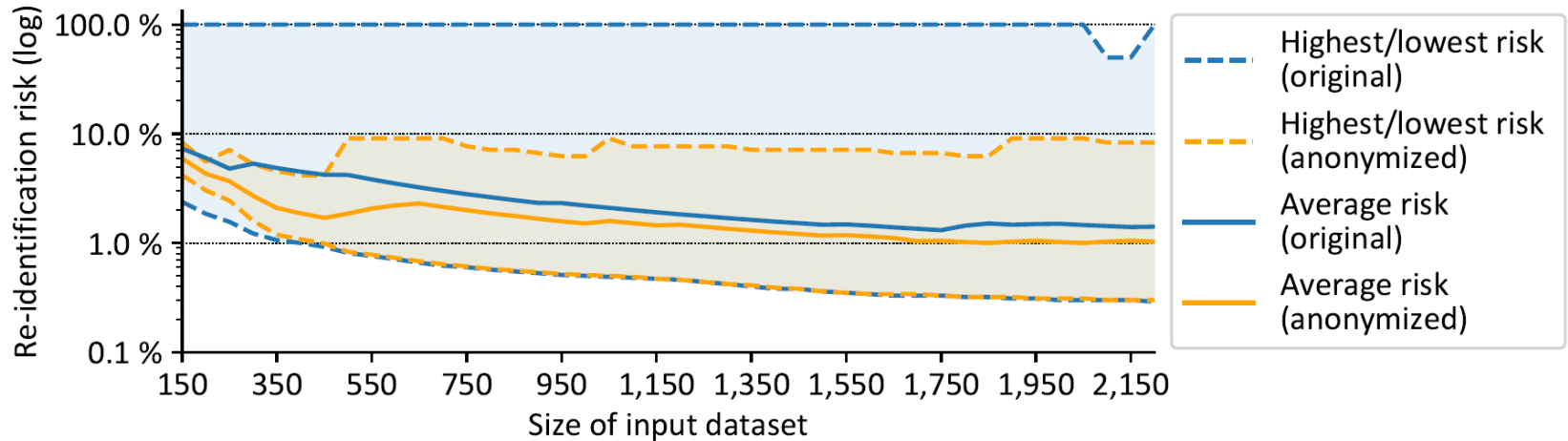
# LEOSS: Evaluation (2)

## Example: descriptive statistics



# LEOSS: Evaluation (3)

## Development of residual risks



# LEOSS: Summary

**Eight additional pipeline stages implement transformations for various modules of the Scientific Use File. Examples:**

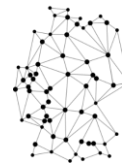
- Categorizing metric variables.
- Making timestamps relative.
- Grouping or suppressing sensitive variables.
  - Modules and stages can be activated dynamically to adjust to needs of different scientific / medical domains.

## **Overall approach**

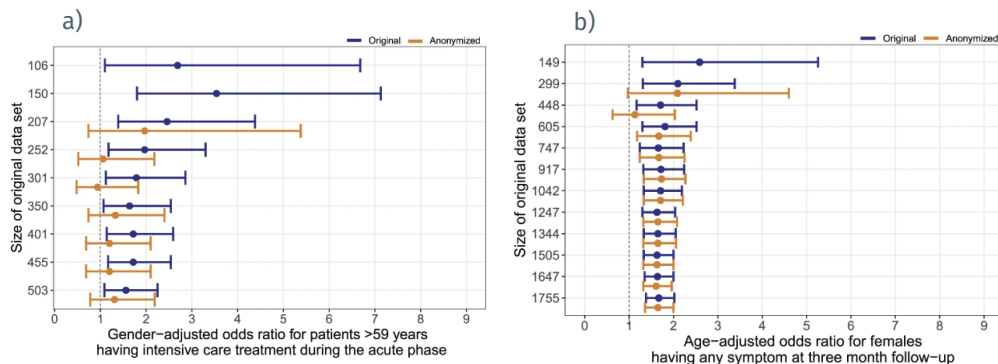
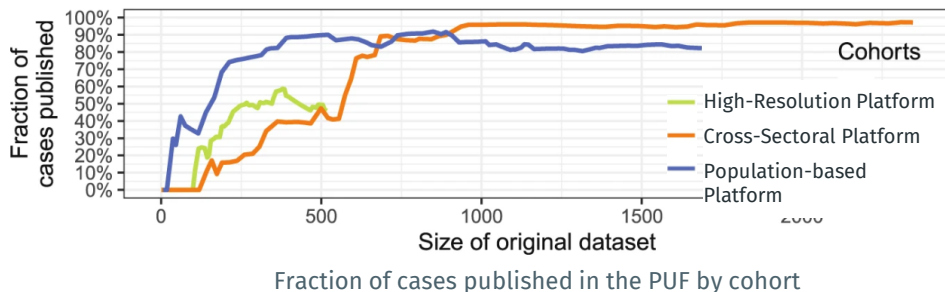
- Context-specific: adopted to the concrete dataset.
- Multiple layers of safeguards: qualitative + quantitative methods.
- Reliance on recommendations from laws and guidelines.
- Risk-based approach requires thorough documentation.



# National pandemic cohort network



**NAPKON**  
NATIONALES  
PANDEMIE  
KOHORTEN  
NETZ



Odds ratios (OR) and 95%-confidence intervals (CI) of patient characteristics in the a) High-Resolution Platform and b) Cross-Sectoral Platform

- **Objective:** Evaluation of how anonymization affects the statistical properties of a Public Use File (PUF)
- **Dataset and risk assessment:** The PUF contains 6.000 cases from 3 different cohorts; It features 16 attributes; 5 of these attributes were associated with an increased risk of re-identification; 8 of the remaining attributes were considered sensitive information
- **Anonymization:** Static generalization of the information; Suppression of records which do not meet the privacy criteria of k-Anonymity with  $k=11$  and t-Closeness with ( $t = 0.5$ )

**Thank you for your attention!**  
**Questions?**