

## Outline



- Motivation & Problem
- Foundations
- 3) Research Questions & Approach
- Results
  - 1) Overview & Classification of De-Identification Methods
  - 2) Requirements
  - 3) Concept Development
- **Conclusion & Limitations**

### Motivation & Problem Statement



#### **Big Data (Analytics)**

- Basis for emerging technologies
- Enormous growth in data
- → Analytics + Data-driven decision-making





#### **Privacy**

- Increased privacy awareness
- Privacy is valued (67 % → no control)<sup>1</sup>
- Maturing privacy laws (GDPR, HIPAA, BDSG)



#### **Wrist-worn Wearable Data**

- Sensitive Health data
- Increasing popularity (36% in Germany)<sup>2</sup>
- Lots of data points
- E.g. heart rates, blood oxygen saturation, location data



<sup>&</sup>lt;sup>1</sup> (European Commission, 2018)

e/anteil-der-smartwatch-nutzer-in-deutschland/

### Motivation & Problem Statement



#### **Risks & Problems:**

- Platform-providers (trustful?)
- 3<sup>rd</sup> party providers (usage / sale)
- no information about storage, processing of data (despite privacy policy)
  - → Disclosure of personal health information
    - → Disclosure of your location data







#### **Wrist-worn Wearable Data**

- Sensitive Health data
- Increasing popularity (36% in Germany)<sup>2</sup>
- Lots of data points
- E.g. heart rates, blood oxygen saturation, location data



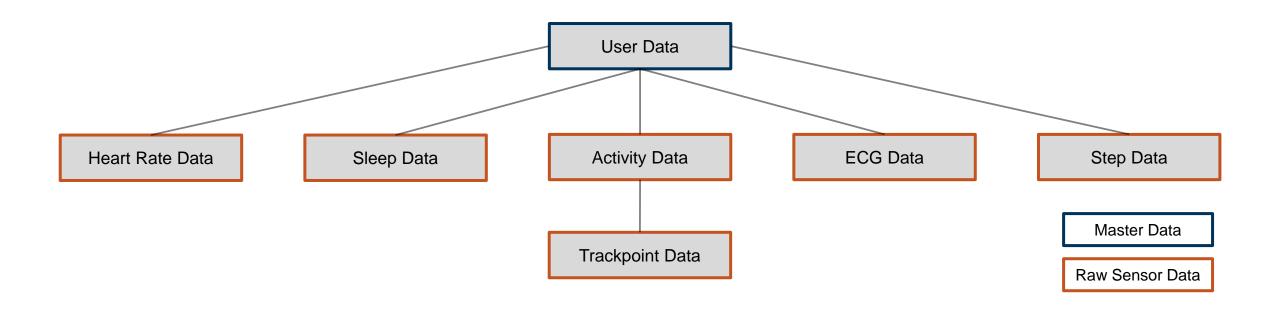
<sup>&</sup>lt;sup>1</sup> (European Commission, 2018)

https://de-statistacom.eaccess.ub.tum.de/statistik/daten/studie/1047586/umfra e/anteil-der-smartwatch-nutzer-in-deutschland/

## Use Case: Wrist-Worn Wearable Data







### **Foundations**



"Privacy "[...] is defined not by what it is, but by what it is not - it is the absence of a privacy breach that defines a state of privacy." (Wu, 2012)

**De-Identification method**: "method for transforming a *dataset* with the objective of reducing the extent to which information is able to be associated with individual *data principals*" (ISO, 2018)

I Terminology of De-Identification: (Tomashchuk et. al., 2019)

→ "a concept of higher level, which covers both anonymization and pseudonymization"

Master Thesis - Final Presentation | Kevin Baumer

## Research Questions / Research Approach



RQ1

What is the **state of the art** of approaches using de-identification methods for privacy-enhancing Big Data Analytics and how can they be distinguished from other approaches?



Status Quo

Literature Research

RQ 2

What are **requirements** for privacy-enhancing analytics of wrist-worn wearable data in the cloud?



Requirements

Literature Research

**Expert interviews** 

Regulations

RQ3

What are **concepts** enabling data privacy for wrist-worn wearable data in the cloud based on **de-identification methods**?



Concepts

Results RQ1 / RQ2

Validation with Experts

Master Thesis - Final Presentation | Kevin Baumer

Research Question 1: What is the state of the art of approaches using de-identification methods for privacyenhancing Big Data Analytics and how can they be distinguished from other approaches?



#### De-Identification – Extensive Literature Review

#### **Databases:**

ScienceDirect







**Search:** methods for De-Identification, Anonymization, Pseudonymization

Total 38 sources



**Comprehensive overviews** 7 sources

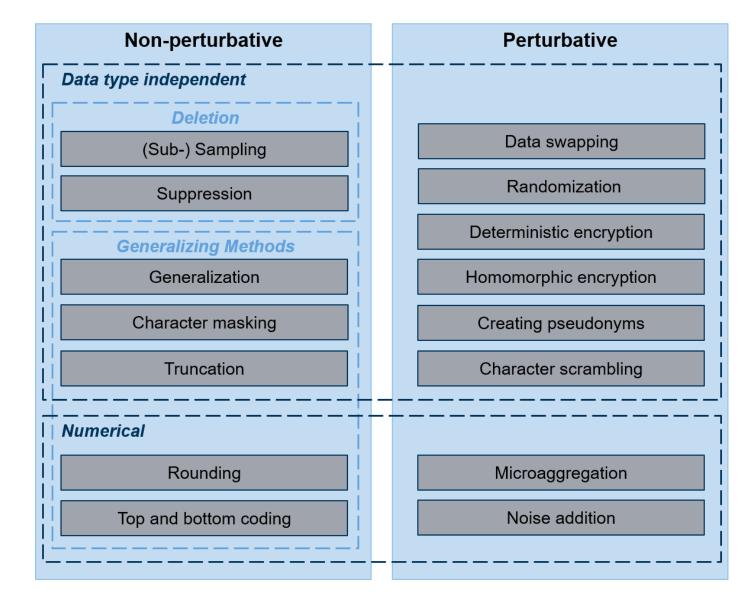
Source	# of methods
Domingo-Ferrer et. al (2019)	6
Mansfield-Devine (2014)	6
Nelson (2015)	14
Tomashchuk (2019)	7
DP Working Party Art. 29 (2014)	5
Bourka, Drogkaris (2018)	5
ISO/IEC 20889 (2018)	16

## Classification of De-Identification Methods



#### Non-perturbative

- → Data stays truthful; Accuracy might be reduced
- → data truthfulness at the record level



#### Perturbative

- → Transformed values not truthful in general
- → statistical properties may be preserved

# Research Question 2: What are requirements for privacy-enhancing analytics of wrist-worn wearable data in the cloud?



## **Expert Interviews**

- 12 semi-structured interviews with industry experts
- Goal: practical insights and implications to derive requirements

ID	Role	Relevant experience (in years)	No. of employees	Duration (hh:mm:ss)
I1	Consultant Data Privacy & Information	>5	1-10	01:04:02
	Security			
I2	Consultant Data Security & Data Privacy	13	1-10	00:56:49
I3	Head of Data Privacy	20	10,001-50,000	00:28:23
<b>I4</b>	Managing Director & Lawyer	11	1-10	00:42:57
I5	Researcher Digital Health	2	11-50	00:28:13
I6	Data-driven Development & Data Pri-	8	>100,000	0:33:43
	vacy Expert			
I7	Information Security Officer & Data Pro-	20	10,001-50,000	00:35:51
	tection Officer			
I8	Head of Data Privacy	23	>100,000	00:39:57
I9	Head of Data Privacy	19	>100,000	00:28:58
I10	Consultant Data Privacy	7	51-250	00:34:30
I11	Chief Information Security Officer	8	251-1,000	00:32:55
I12	Key Expert Data Privacy	6	>100,000	00:52:21

## **Identified Requirements**



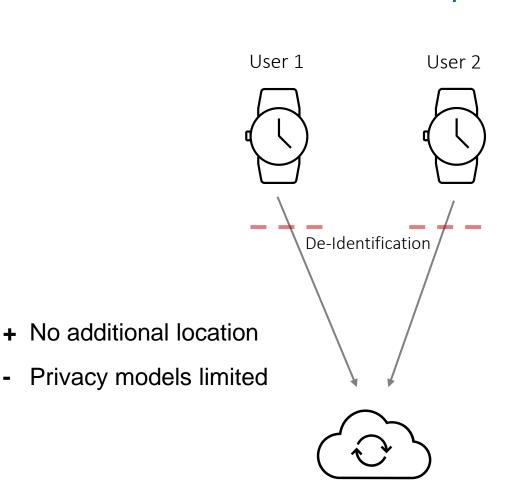
Local transformation of data Low performance overhead k-anonymity enforcement Privacy levels 3 Generic wrist-worn data model 5 Compliance with regulations 10

Constraints from a privacy perspective Protection against complete disclosure Transparency Transformation of identifiers

# Research Question 3: What are concepts enabling data privacy for wrist-worn wearable data in the cloud based on de-identification methods?

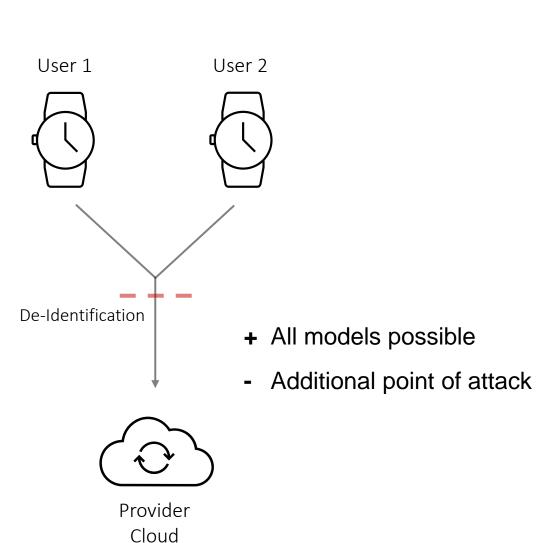


## Technical Architecture – Two Options



Provider

Cloud



## K-anonymity & AVC



**k-anonymity:** A table satisfies k-anonymity if every record in the table is indistinguishable from at least k-1 other records with respect to every set of quasi-identifier attributes. (Sweeney, 2002)

Attribute value combinations (AVC): The attribute value combinations represent the total number of possible value combinations for all attributes in a data set.

Zip Code	Age	Nationality
13053	28	Russian
13068	29	American
13068	21	Japanese
13053	23	American
14853	50	Indian
14853	55	Russian
14850	47	American
14850	49	American
13053	31	American
13053	37	Indian
13068	36	Japanese
13068	35	American

Zip Code	Age	Nationali
130**	< 30	*
130**	< 30	*
130**	< 30	*
130**	< 30	*
1485*	$\geq 40$	*
130**	3*	*
130**	3*	*
130**	3*	*
130**	3*	*

<sup>(</sup>a) k-anonymous table with k=1

Attribute	# Distinct values
Zip Code	20
Age	80
Nationality	10

$$\rightarrow$$
 AVC = 20 \* 80 \* 10 = 16,000

→ Concept for local probabilistic k-anonymity

<sup>(</sup>b) k-anonymous table with k=4

## Local probabilistic k-anonymity: Monte Carlo simulation



Monte Carlo simulation: repeated random sampling to obtain numerical results

### 3 parameters:

n: number of iterations

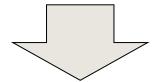
## Local probabilistic k-anonymity: Monte Carlo simulation (Python)



#### Input parameters

$$g = 5$$
 (AVC)

$$r = 200$$
 (records)

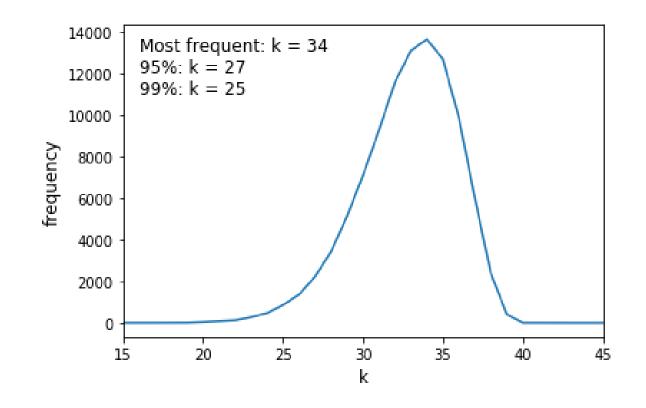


#### <u>Results</u>

99%: k = 25

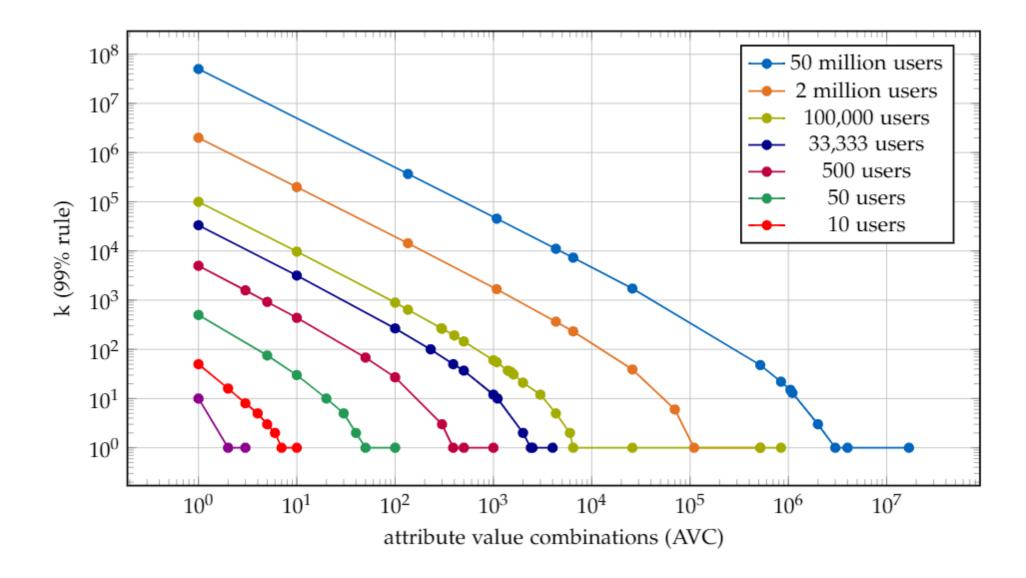
→ 1% risk remains

(anonymity threshold)



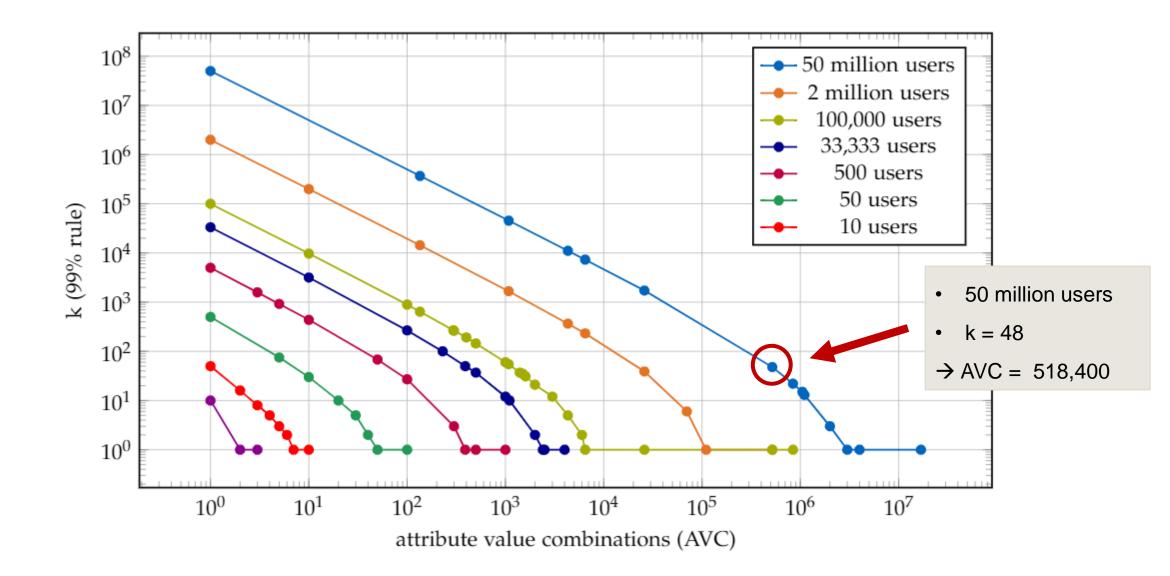
## Local probabilistic k-anonymity - Results





## Local probabilistic k-anonymity - Results





## De-identification methods to reduce AVC



Quasi Identifiers	Distinct values (initially)	
handedness	2	
gender	3	
height	51	
weight	81	
country	195	
currentGear	2,000	
createdDate	3,650	
birthdate	15,695	
eMail	~ ∞	
firstName	~ ∞	
lastName	~ ∞	
profileImage	~ ∞	
AVC	~ ∞	

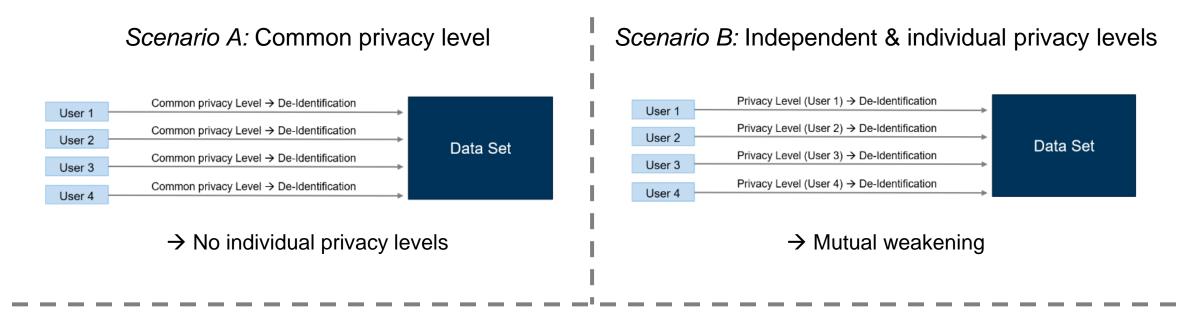
Methods
Suppression
-
Generalization (5 cm)
Generalization (5 kg)
Generalization (continent)
Generalization (brand)
Suppression
Generalization (5 years)
Creating pseudonym
Suppression
Suppression
Suppression

Distinct values
-
3
10
16
6
20
-
9
-
-
-
-
518,400

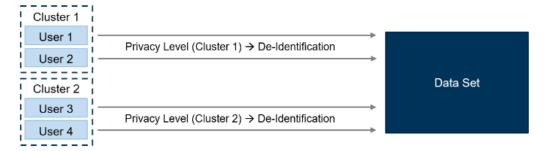
Master Thesis - Final Presentation | Kevin Baumer

## Local probabilistic k-anonymity – 4 Scenarios





#### Scenario C / D: Privacy clusters with equal / unequal distribution



→ Best approach to achieve individual privacy levels

Master Thesis - Final Presentation | Kevin Baumer © sebis 1

## Implications and Results for wrist-worn wearable data



- Proposed approach: privacy clusters
- Limiting the number of different clusters (low / medium / high)
- Optimization problem to investigate combination of clusters (unequal distribution)
- Consider fluctuations between clusters (e.g. risk calculation with 80% of users)

## Summary & Results



RQ 1

#### Complete & comprehensive overview of de-identification methods

- → Clear picture and impression of available methods
- → Relation and differentiation between methods
- → Supporting decision process for method choice

RQ 2

#### 10 identified concept requirements

→ Generic wrist-worn data model

RQ3

### Local probabilistic k-anonymity concept for wrist-worn wearable data

- → Privacy estimates based on Monte Carlo simulation
- → Concept involving privacy clusters

## Limitations & Future Work



### Limitations

- Lacking availability of a data set
- Validation in real application scenario
- Suitable domain experts

## Future work

- Extension of the local probabilistic k-anonymity concept
- Evaluation and testing on data set
- Benchmark against local differential privacy

