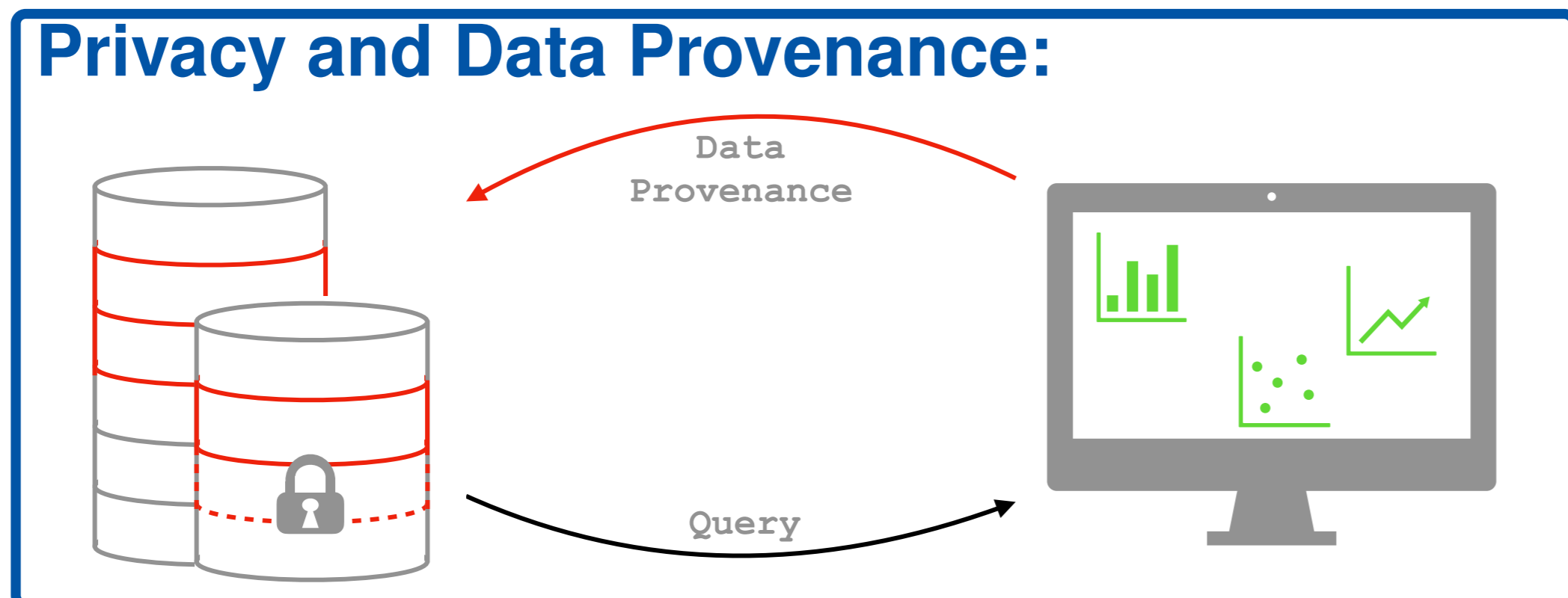


Privacy Aspects of Provenance Queries

Motivation



Privacy:

- protection of personal data against unauthorized collection, storage and publication
- possibility of not re-identifying single persons in a bunch of data

- big data: amount of data \uparrow , transparency \downarrow
- GDPR: data protection more important than ever before

Data Provenance:

- lineage of data
- **where**: Where does the data come from?
 \Rightarrow names of the source relations like Grades
- **why**: Why was this result achieved?
 \Rightarrow witness bases like $\{\{t_4, t_5, t_6\}\}$
- **how**: How was the result calculated?
 \Rightarrow provenance polynomials like $\frac{4.0 t_4 + 5.0 t_5 + 1.3 t_6}{t_4 + t_5 + t_6}$

Possible provenance-based Database Reconstructions:

Grades	ID	Module	Grade
t_1	13	Mathematics	1.0 A
t_2	27	Data Science	2.3 B
t_3	27	Theory	1.7 B
t_4	115	Mathematics	4.0 D
t_5	115	Data Science	5.0 E
t_6	115	Law	1.3 A

ID	Module	Grade
13	η_1	1.0
27	η_2	2.0
115	η_3	3.3

where

ID	Module	Grade
13	η_1	1.0
27	η_2	2.0
27	η_3	2.0
115	η_4	3.3
115	η_5	3.3
115	η_6	3.3

why

ID	Module	Grade
13	η_1	1.0 A
27	η_2	2.3 B
27	η_3	1.7 B
115	η_4	4.0 D
115	η_5	5.0 E
115	η_6	1.3 A

how

ID	avgGrade
13	1.0 A
27	2.0 B
115	3.3 C

where Grades
why $\{\{t_4, t_5, t_6\}\}$
how $\frac{4.0 t_4 + 5.0 t_5 + 1.3 t_6}{t_4 + t_5 + t_6}$

Data Protection Problems with *where*, *why* and *how*

- **where**: (1) no data worth protecting available or (2) save the tuple itself \Rightarrow privacy aspects negligible or a huge problem
- **why**: if distribution of data is known and data is equal \Rightarrow privacy aspects could be a problem
- **how**: too much information recoverable \Rightarrow privacy aspects are in all probability a problem

Solution Approaches:

No. 1: Generalization

No. 2: Intensional answers

No. 3: Differential privacy

No. 4: Permutation

$$\frac{4.0 \times t_4 + 5.0 \times t_5 + 1.3 \times t_6}{t_4 + t_5 + t_6} = \frac{1.3 \times t_6 + 4.0 \times t_4 + 5.0 \times t_5}{t_6 + t_4 + t_5}$$