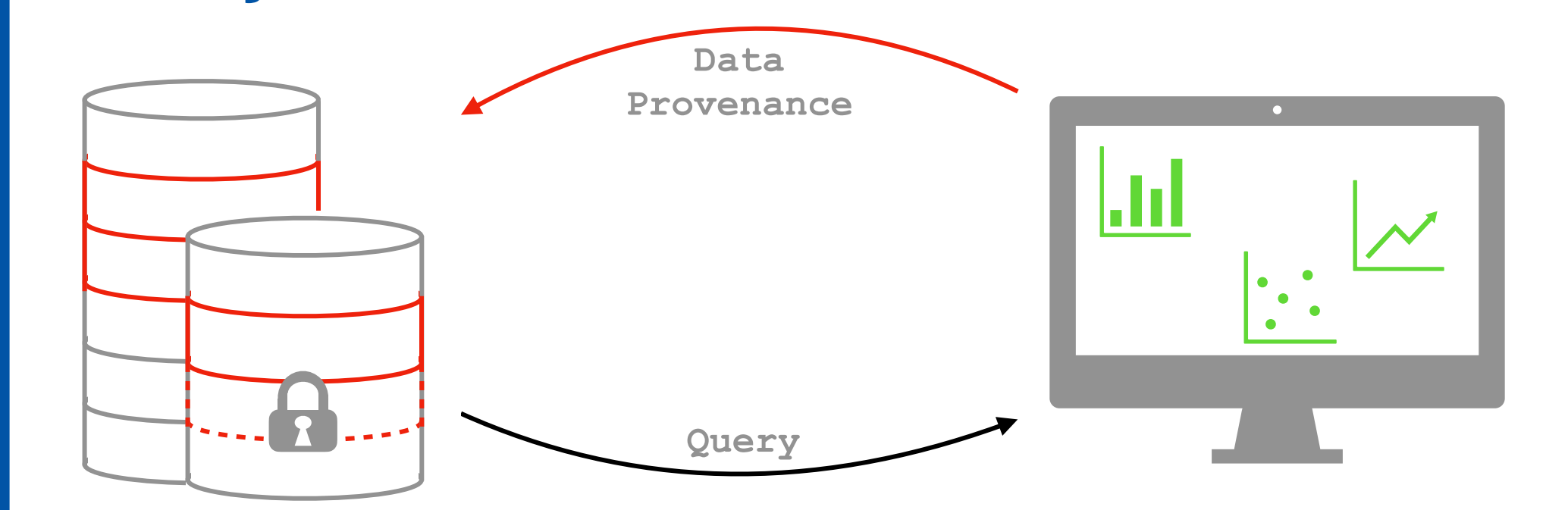


Privacy Aspects of Provenance Queries

Motivation

Privacy and Data Provenance:



Privacy:

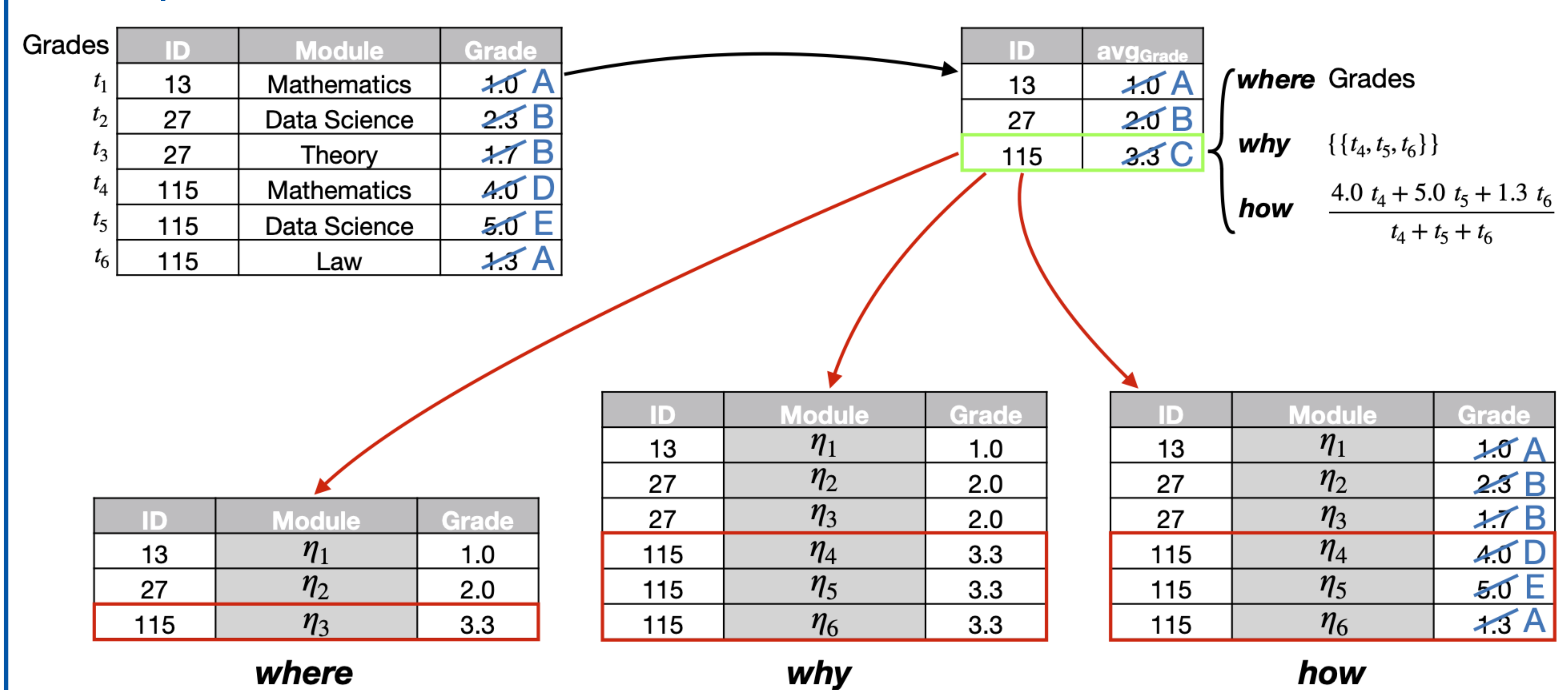
- protection of personal data against unauthorized collection, storage and publication
- possibility of not re-identifying single persons in a bunch of data

- 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:



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:

