

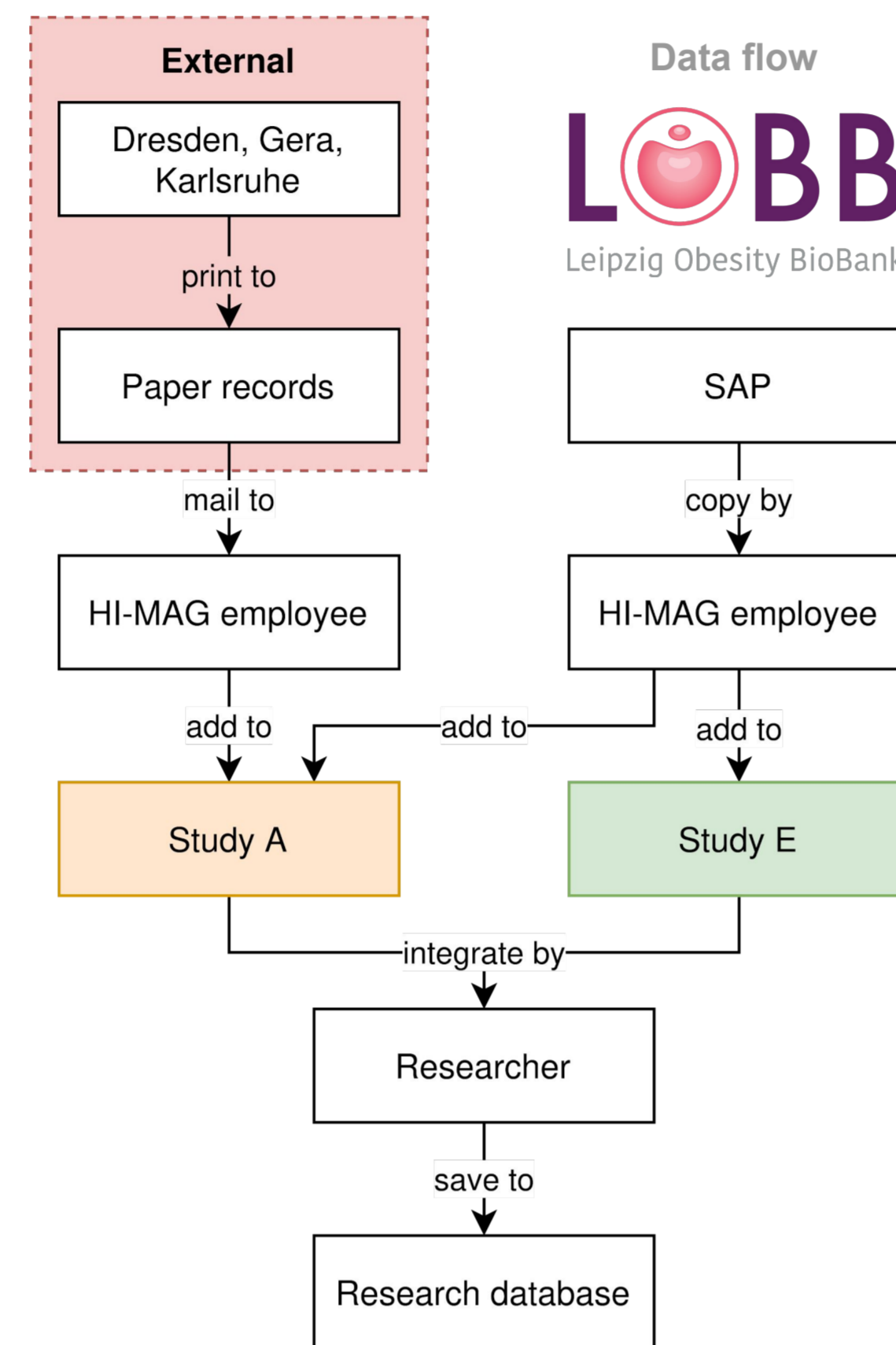
Improving record linkage quality on identification data in the Leipzig Obesity BioBank

Maximilian Jugl^{1,2}, Anne Hoffmann³, Matthias Kern^{3,4}, Matthias Blüher^{3,4}, Thomas Ebert⁴, Toralf Kirsten^{1,2}

¹University of Leipzig Medical Center, Medical Informatics Center, Dept. Medical Data Science, Leipzig, ²Leipzig University, Institute of Medical Informatics, Statistics, and Epidemiology, Leipzig, Germany, ³Helmholtz Institute for Metabolic, Obesity and Vascular Research (HI-MAG), Helmholtz Zentrum München, University of Leipzig and University Hospital Leipzig, Leipzig, Germany, ⁴University of Leipzig Medical Center, Department of Endocrinology, Nephrology and Rheumatology, Leipzig, Germany

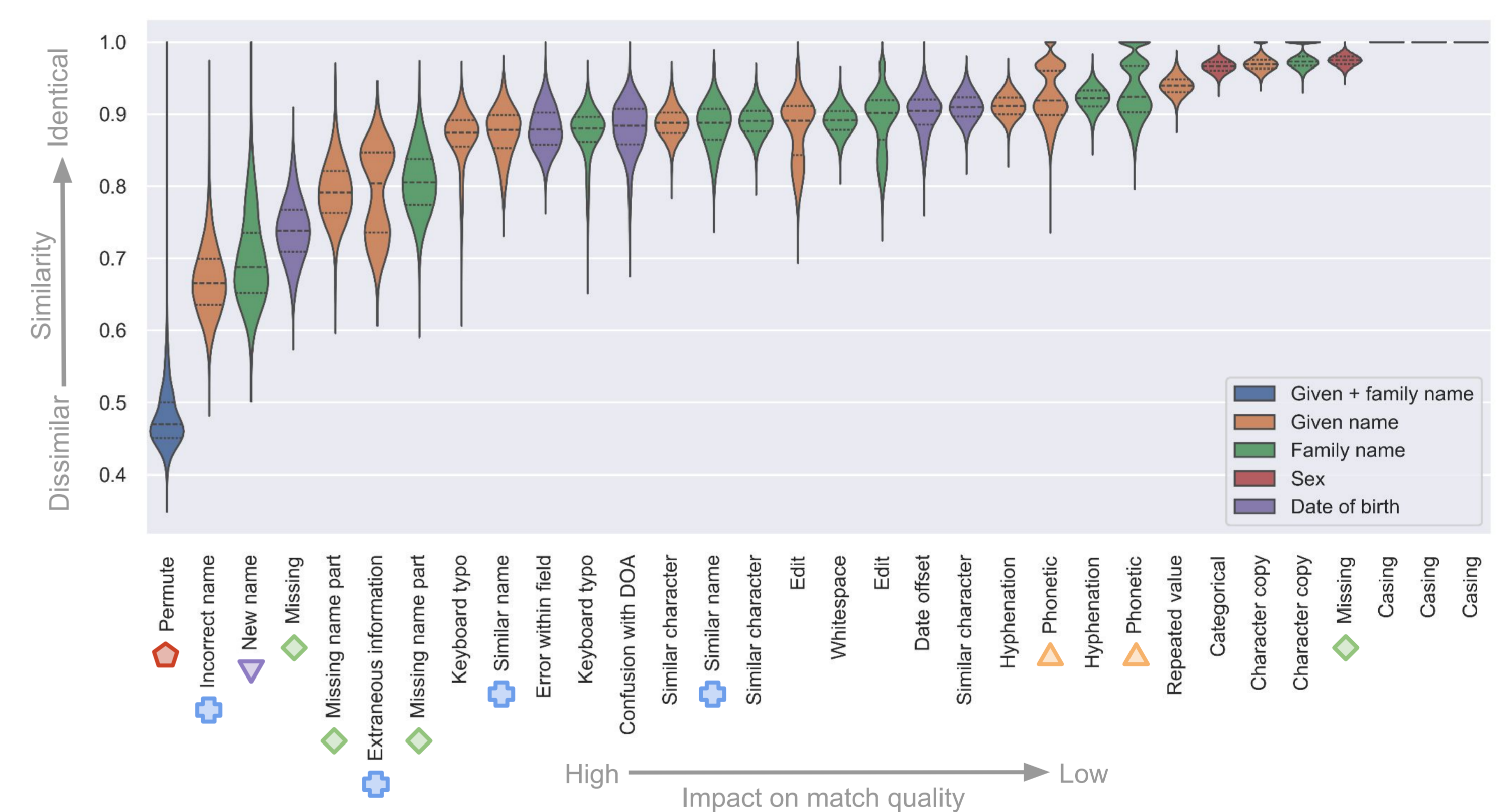
Background

- Leipzig Obesity BioBank (LOBB): longitudinal medical study researching diseases related to obesity with 8,000+ patients
- LOBB seeks to integrate medical data from multiple sources with varying input methods
- Privacy-preserving record linkage (PPRL) requires real-world data to test algorithms¹
- Few studies based on real-world identification data (IDAT)^{2,3,4} ⇒ generation of realistic data feasible⁵
- Use LOBB to analyse error sources and generate realistic data generation workflow**



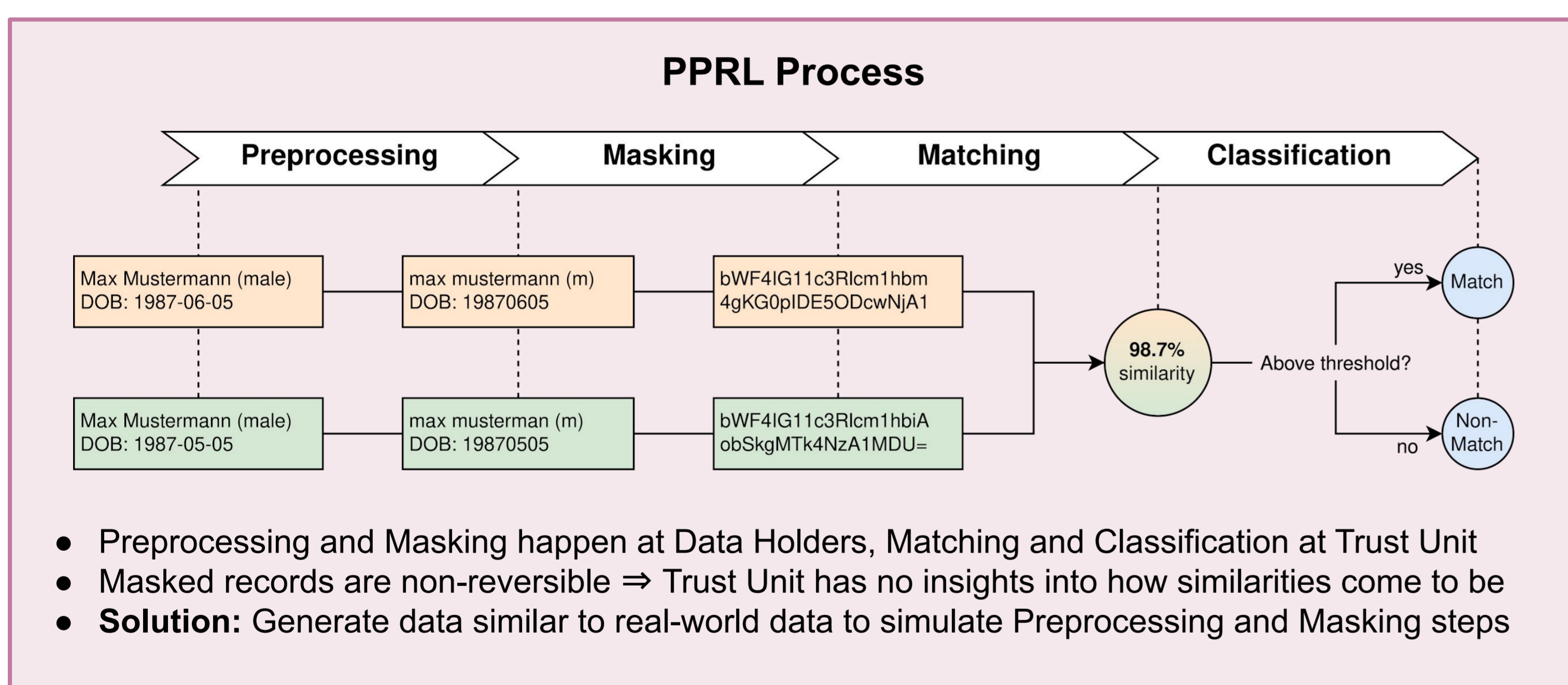
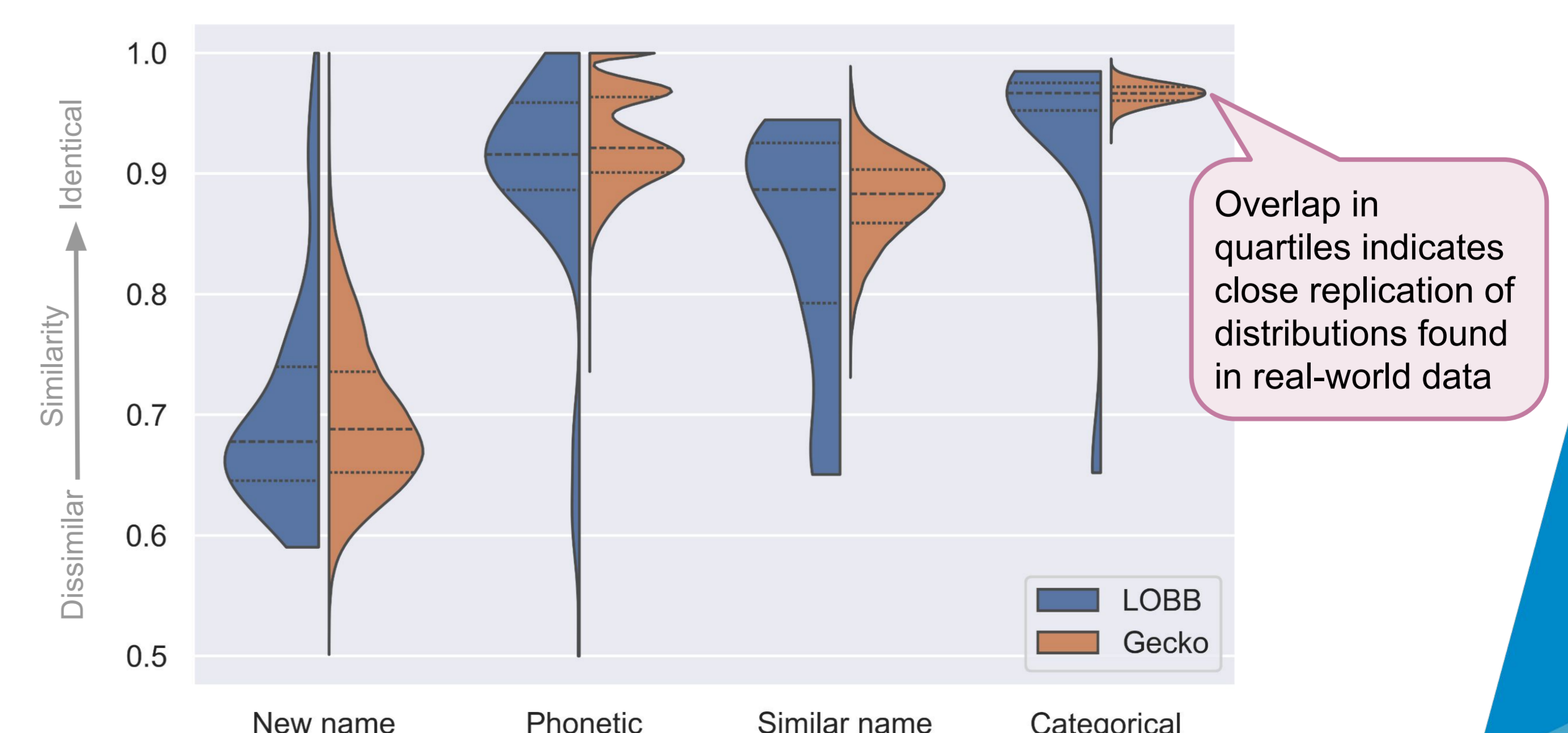
Results

- 275 erroneous records in total ⇒ **1.94% error rate**
- 22 error classes with phonetic errors[▲] and new family names[▼] taking up **over 48% of observed errors**
- Erroneous[⊕] or missing values[◇] and permutation of given and last names[◊] rarely occur, but have a **drastic impact on match quality**
- Gecko*-generated data closely replicates LOBB data in computed similarities of clean to mutated records



Conclusion

- Reproducible workflow for generating realistic data in a medical study from multiple separate data sources
- Overview of error classes from an authentic source
- Proof of *Gecko*'s capabilities to generate and mutate realistic data found in the real world



Methods

- Determine sources of typographic errors between personal record pairs in two LOBB sub-studies
- Derive error configuration for realistic data generation tool *Gecko*⁶, then generate and mutate data
- Compare computed similarities between record pairs in LOBB and those generated by *Gecko*