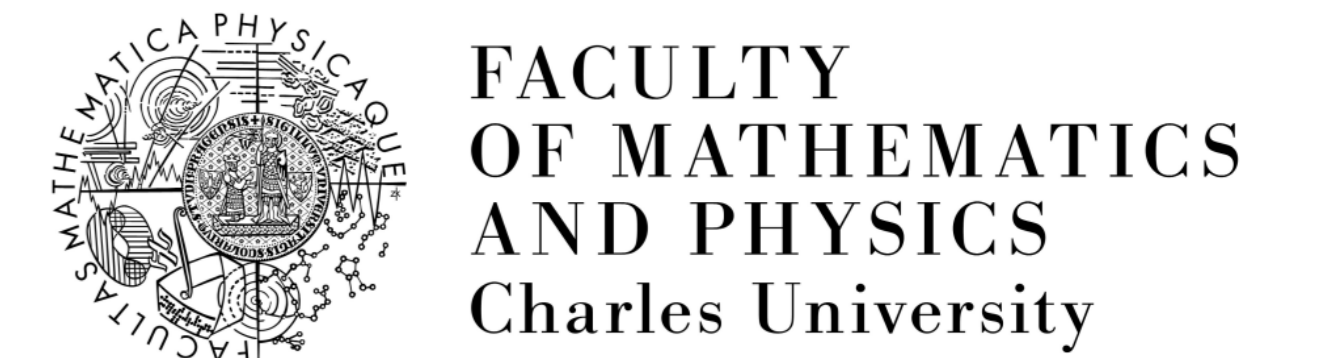
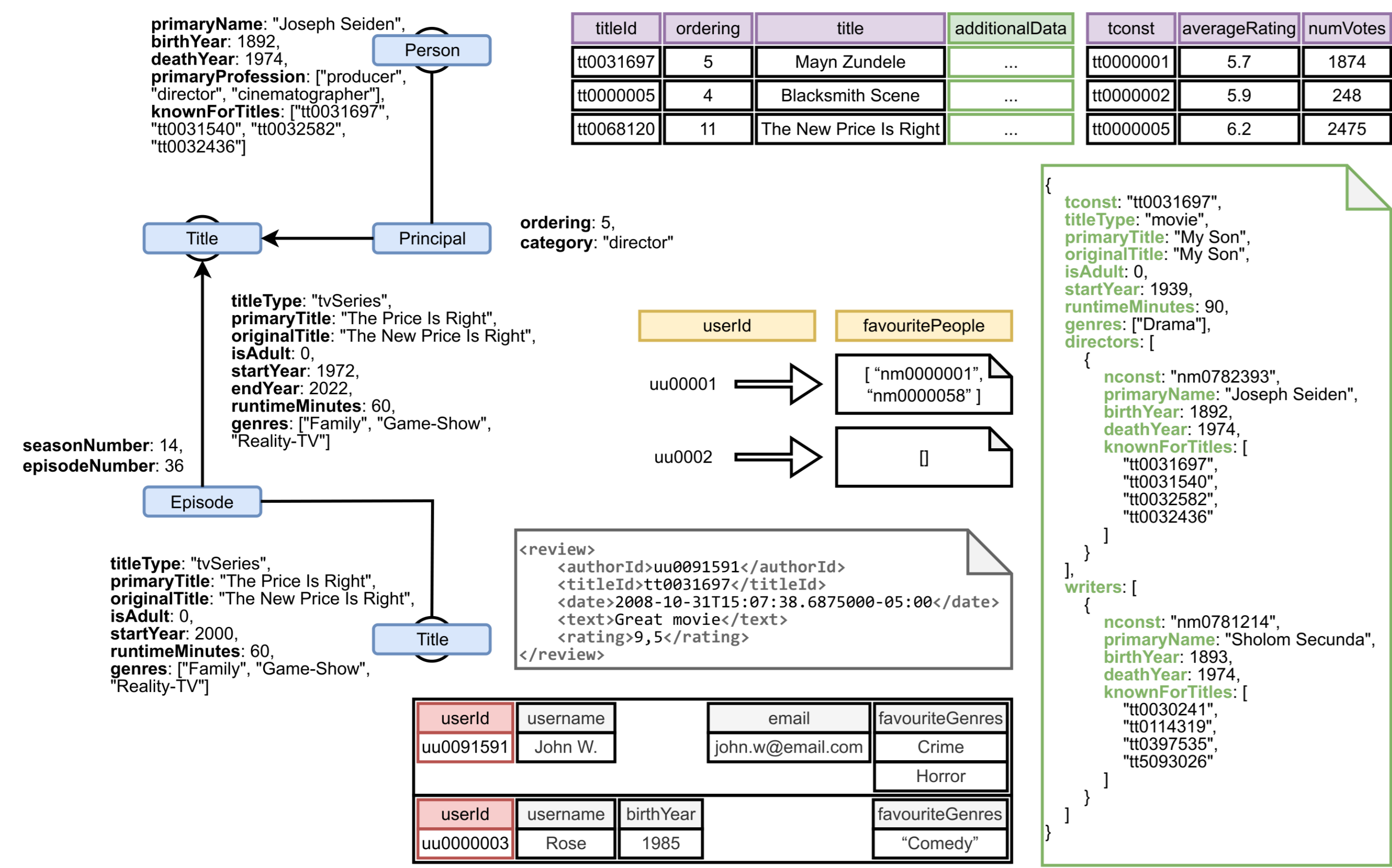


SCHEMA INFERENCE FOR MULTI-MODEL DATA

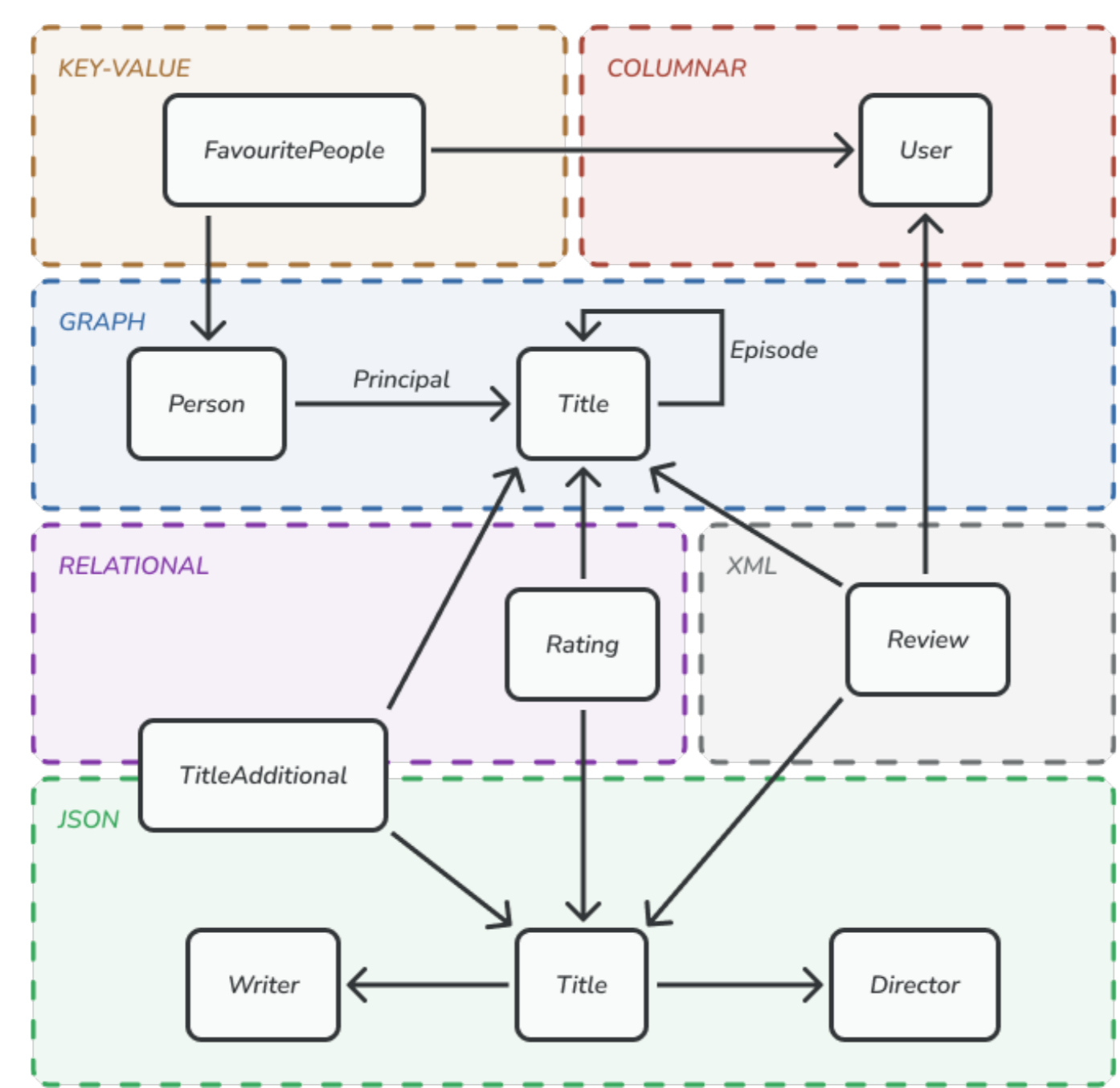
Author: Mgr. Sebastián Hricko, Supervisor: doc. RNDr. Irena Holubová, Ph.D., Advisor: Ing. Pavel Koupil
 Department of Software Engineering, Faculty of Mathematics and Physics, Charles University



PROBLEM DEFINITION

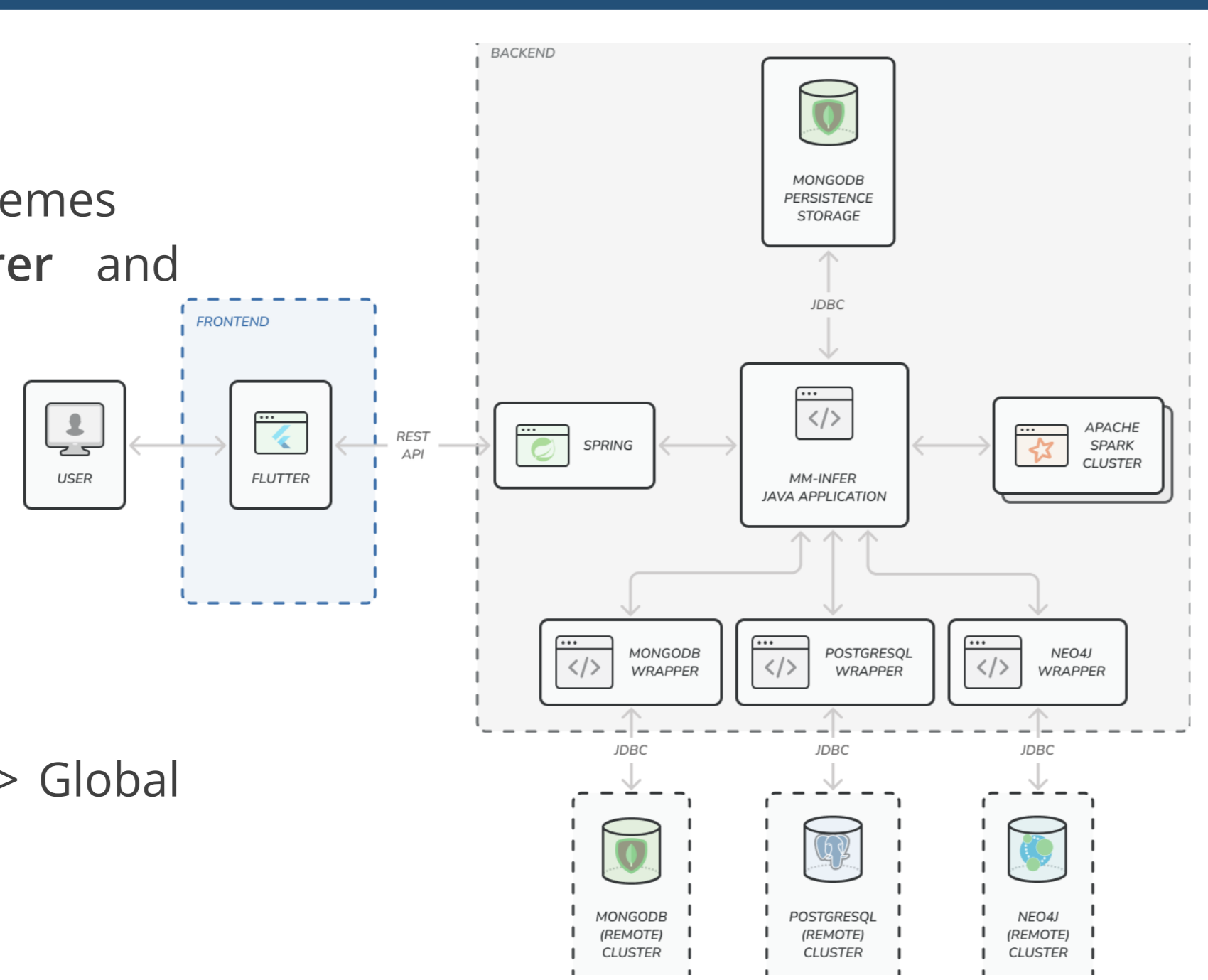


- Consider a multi-model scenario, where the data storage model corresponds to the expected use of the data (Document, Key-Value, Columnar, Relational, Graph, ...)
- Questions we aim to answer:
 - What is the schema of the data in each model?
 - How is the data interconnected between models?
- We want to obtain a global view of the data model and identify references and redundancies in the dataset



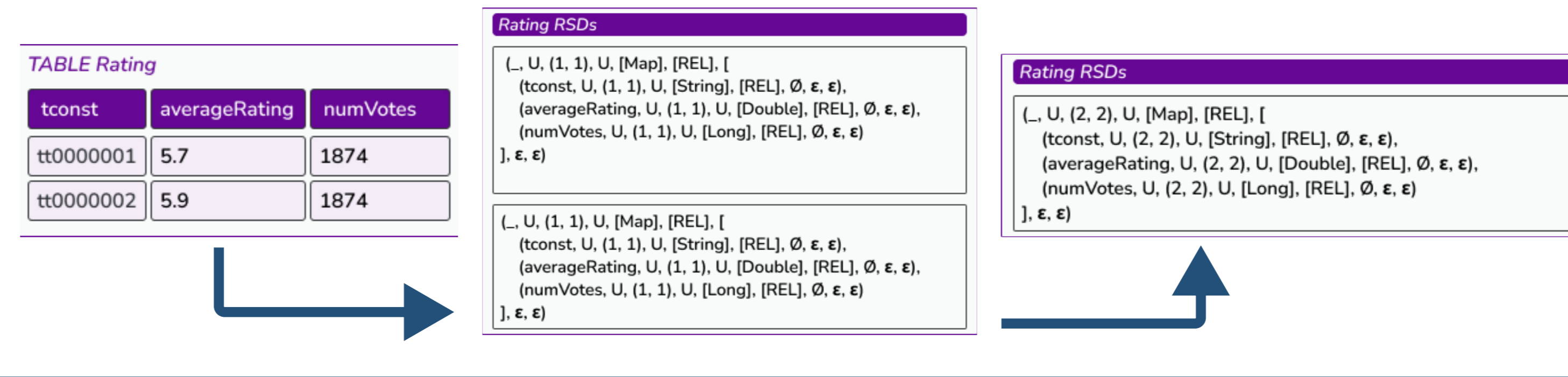
PROPOSED APPROACH

- Unification of models
- Analysis and reduction of unified schemes
 - Algorithms Record-based Inferer and Property-based Inferer
 - Apache Spark MapReduce
- Discovering candidates
 - IO, references, redundancies
 - Candidate Miner algorithm
 - Heuristics
- Evaluation of candidates
- Connecting reduced local schemes -> Global scheme



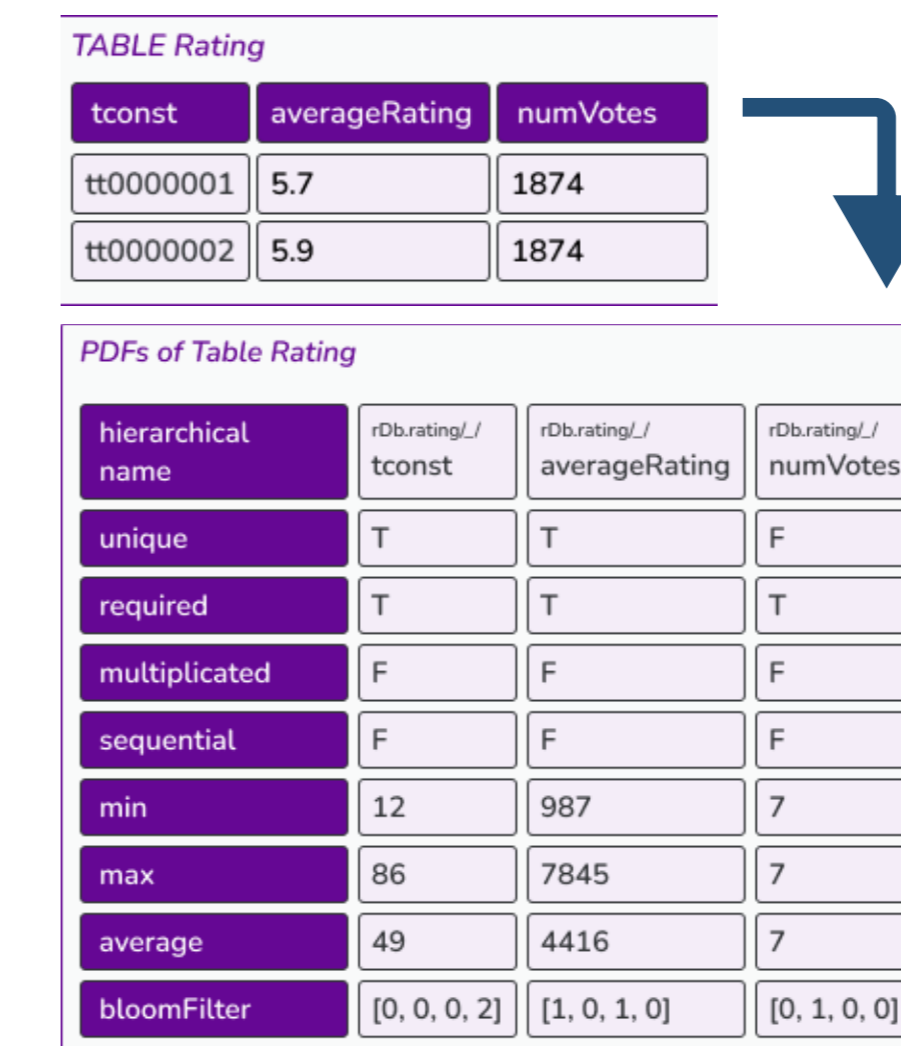
RECORD-BASED INFERRER

- Traversing collections by records (i.e. PostgreSQL table - rows)
- For each record, a simple scheme (RSD) copying the structure (recursively) is created
- Schemas of individual records are merged into the overall description of the collection



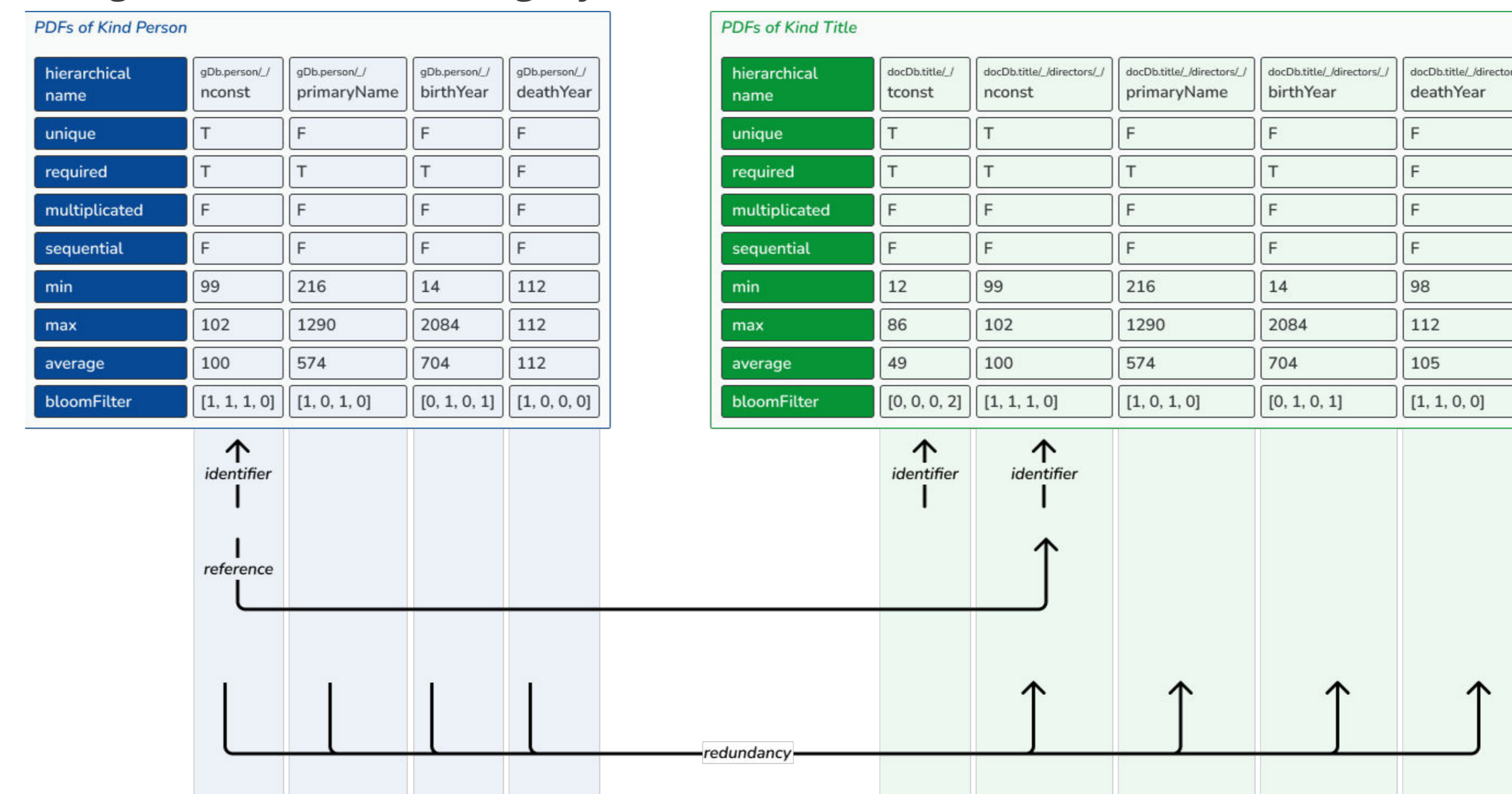
PROPERTY-BASED INFERRER

- Traversing collections by attributes/properties
 - (PostgreSQL table - columns in rows)
- For each property, a property domain footprint (PDF) is created
- A structure-ignoring RSD is generated for each property (without recursion)
- A schema of the collection is created based on the hierarchical names and PDFs



CANDIDATE MINER

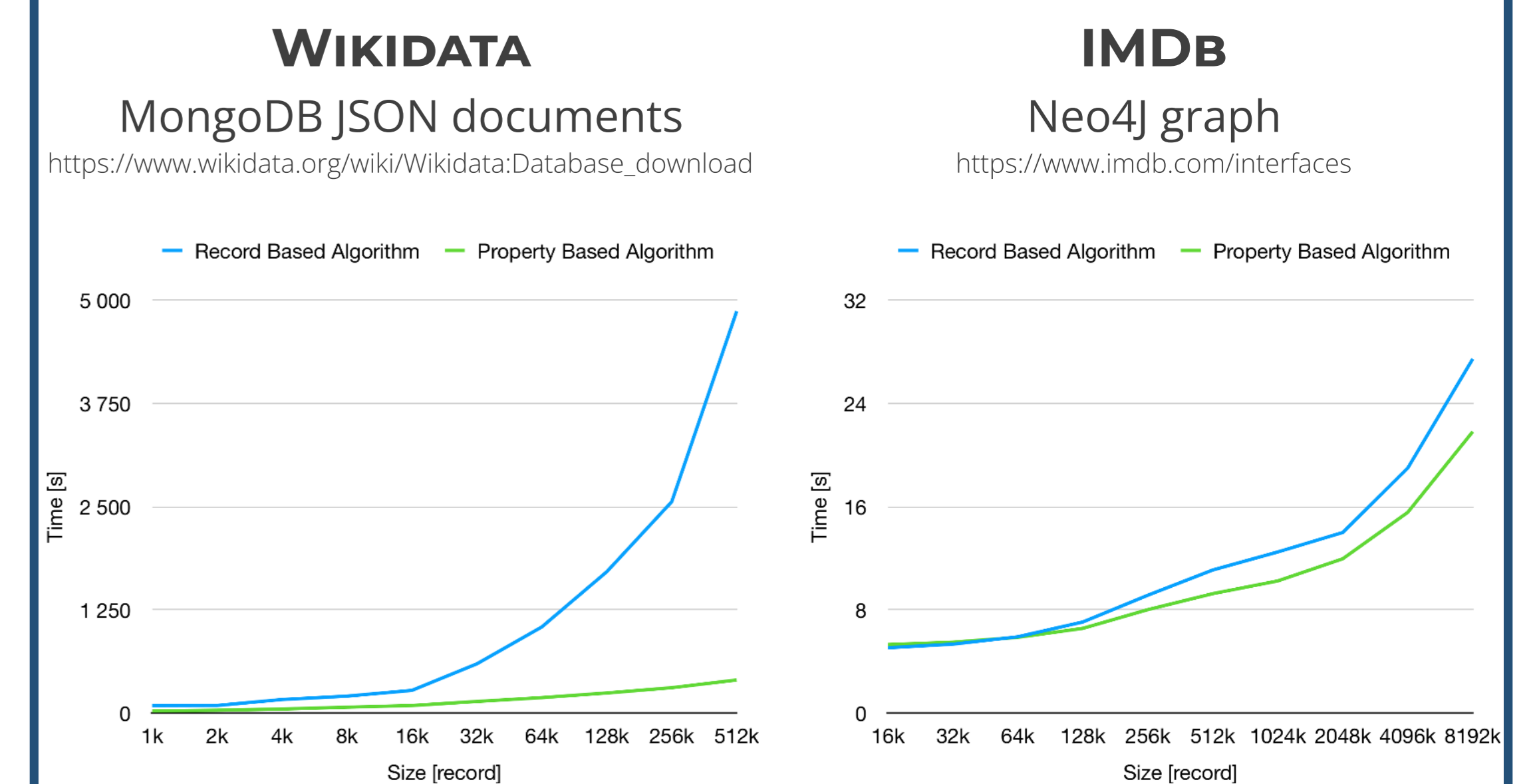
- Based on the Property-based Inferer
- We are looking for candidates for unique identifiers
- We compare the generated PDFs to get references
- We look for overlaps in values of size at least k to reveal redundancies
- We get candidates for Integrity Constraints



CONTRIBUTIONS

- First approach for multi-model schema inference
 - Process the data in an agnostic way of the underlying model
 - It can be used to generate data schemas in commonly used models (relational, document (XML, JSON), key-value, columnar, graph)
 - Finding interconnections between data within models and between models
 - Revealing redundancy in data
- Modular = easily extensible
- Scalable, supporting distributed computing
- Extends schema generation possibilities for single-model data (uniqueness, regular expressions, more advanced references and redundancies)
- The proposed algorithms were experimentally verified

EXPERIMENTAL EVALUATION



PUBLICATION ACTIVITY

- The results of the work were published (accepted) within conferences:
 - EDBT: 25th International Conference on Extending Database Technology
 - ACM / IEEE 25th International Conference on Model Driven Engineering Languages and Systems (MODELS)
- Pavel Koupil, Sebastián Hricko, and Irena Holubová. Mm-infer: A tool for inference of multi-model schemas. In Julia Stoyanovich, Jens Teubner, Paolo Guagliardo, Milos Nikolic, Andreas Pieris, Jan Mühlig, Fatma Özcan, Sebastian Schelter, H. V. Jagadish, and Meihui Zhang, editors, Proceedings of the 25th International Conference on Extending Database Technology, EDBT 2022, Edinburgh, UK, March 29 - April 1, 2022, pages 2:566-2:569. OpenProceedings.org, 2022
- Pavel Koupil, Sebastián Hricko, and Irena Holubová. A universal approach for multi-model schema inference. J. Big Data, 9(1):97, 2022.61