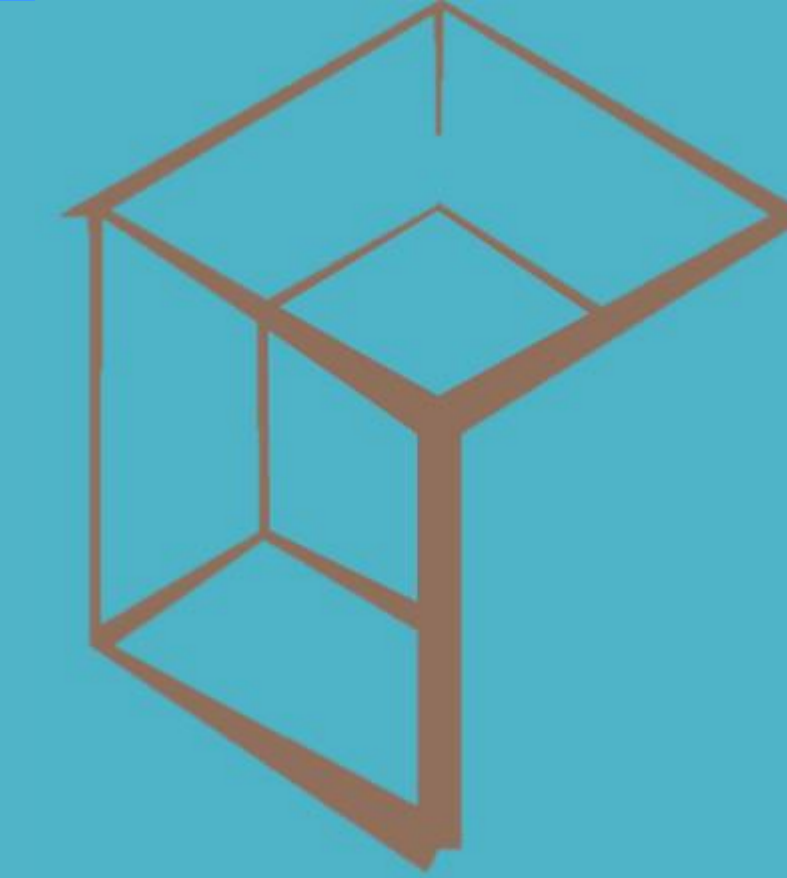


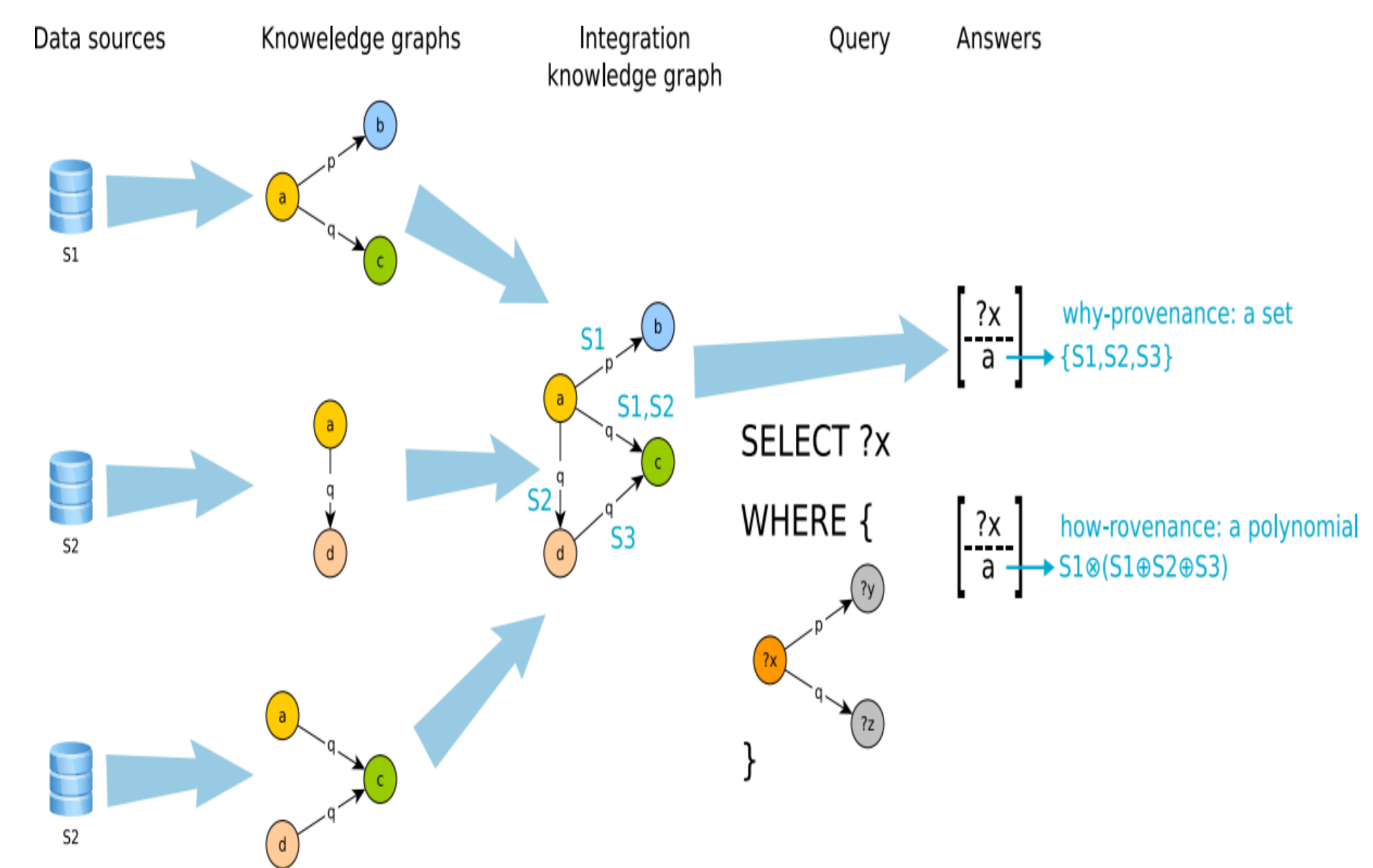
# NPCS: Native Provenance Computation for SPARQL



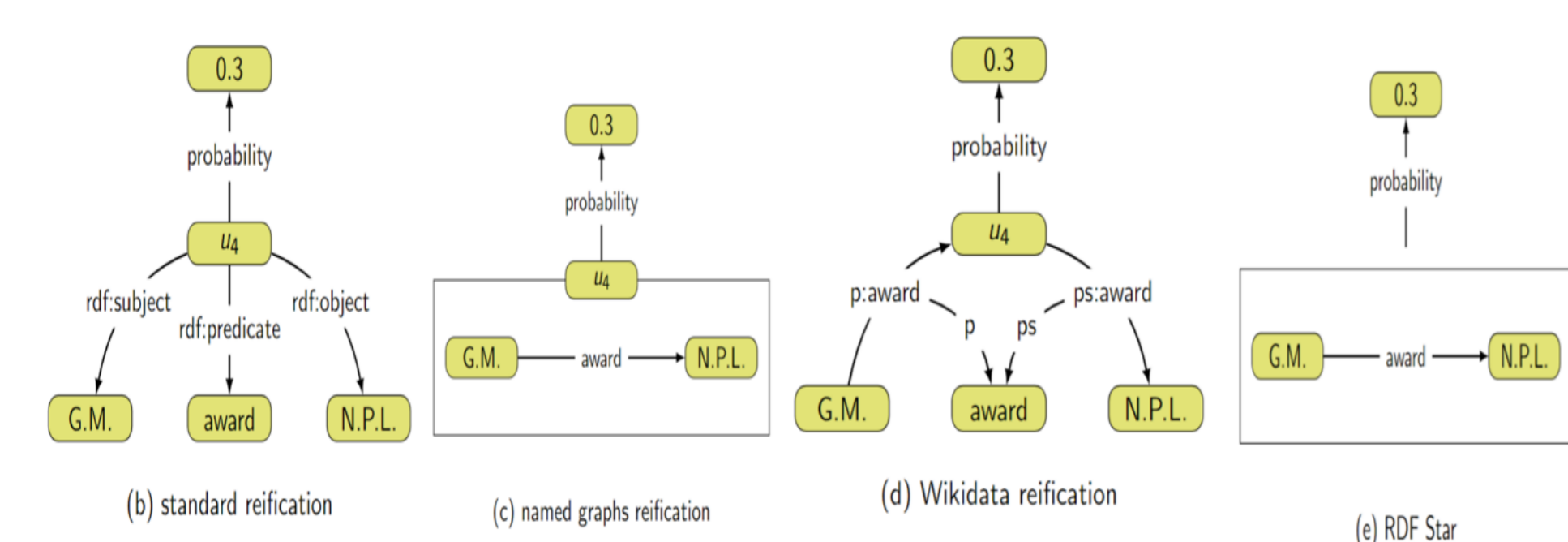
Zubaria Asma<sup>1,2</sup>, Daniel Hernández<sup>3</sup>, Luis Galárraga<sup>4</sup>, Giorgos Flouris<sup>1</sup>, Irini Fundulaki<sup>1</sup> and Katja Hose<sup>5</sup>  
 {<sup>1</sup>FORTH-ICS,<sup>2</sup>University of Crete} Heraklion, Crete, Greece, <sup>3</sup>University of Stuttgart, Stuttgart, Germany, <sup>4</sup>Inria, Rennes, France, <sup>5</sup>TU Wien, Wien, Austria

## Introduction

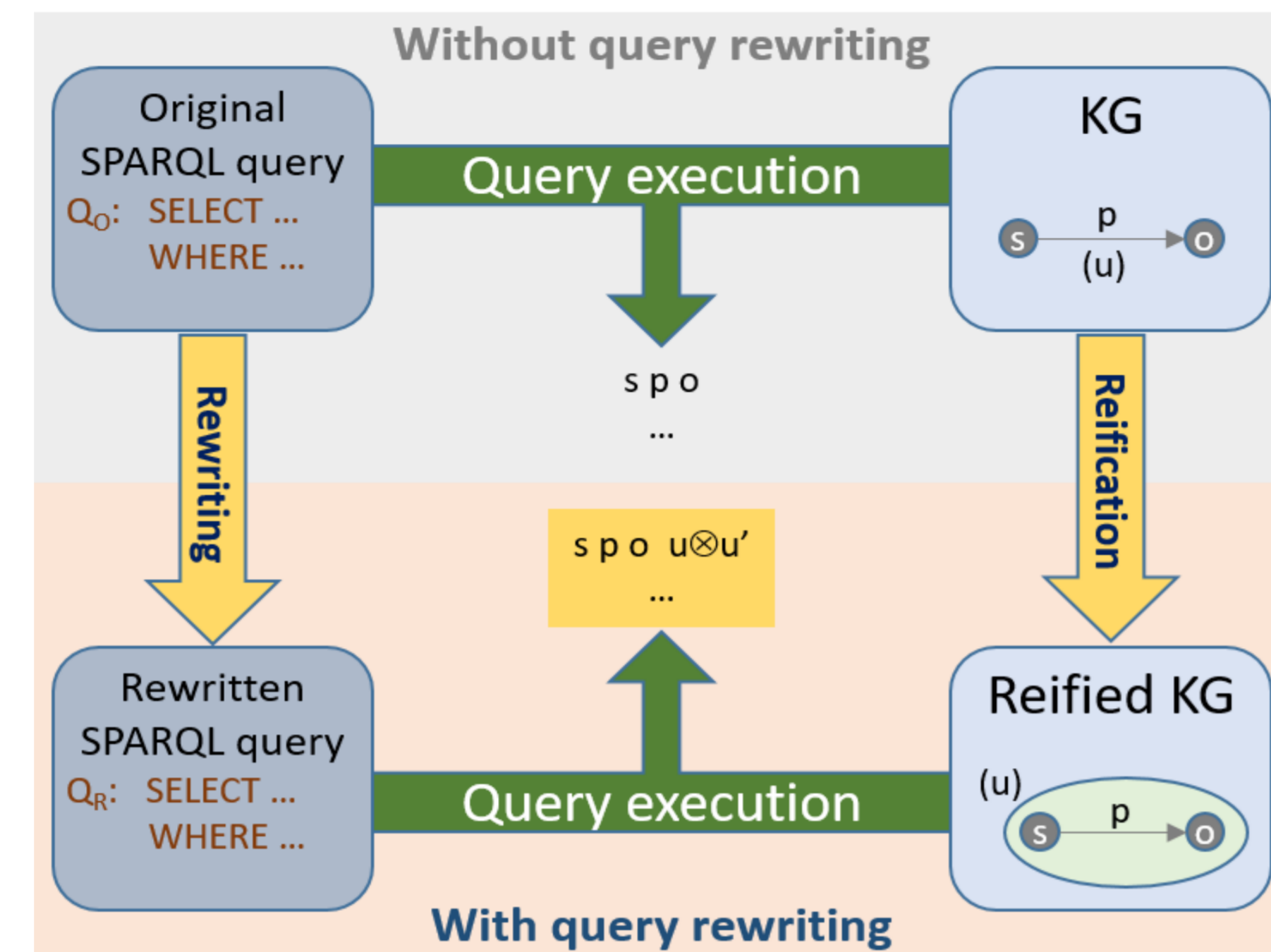
- ❖ Knowledge Graphs (KGs)
  - Important in both academia and industry
  - Ideal for integrating data from various sources
- ❖ Provenance
  - Critical for trust assessment and dynamic data
- ❖ Our Contribution (NPCS)
  - Enriches results with how-provenance annotations
  - Supports monotonic and non-monotonic SPARQL
  - NPCS rewrites SPARQL query  $Q$  into  $Q'$ , generating how-provenance polynomials in the spm-semiring  $K = (K, \oplus, \otimes, \ominus, 0, 1)$  for K-graph  $G$ .



## NPCS supports these reification schemes



## Schema for NPCS architecture



**Query: Get all women awarded with nobel prize in Literature**

**SPARQL Query**  
 SELECT ?x  
 WHERE {  
 ?x gender female .  
 ?x award "N.P.Literature"  
 }

**Result**  
 [ ?x  
 G. Mistral  
 O. Tokarczuk ]

## NPCS Query

(SELECT ?x AggSum(?prov))  
 AS ?k  
 WHERE {  
 Reify(?x, gender, female, ?prov0).  
 Reify(?x, award, "N.P.Literature", ?prov1)  
 BIND(Prod(?prov0, ?prov1) as ?prov)  
 } group by(?x)

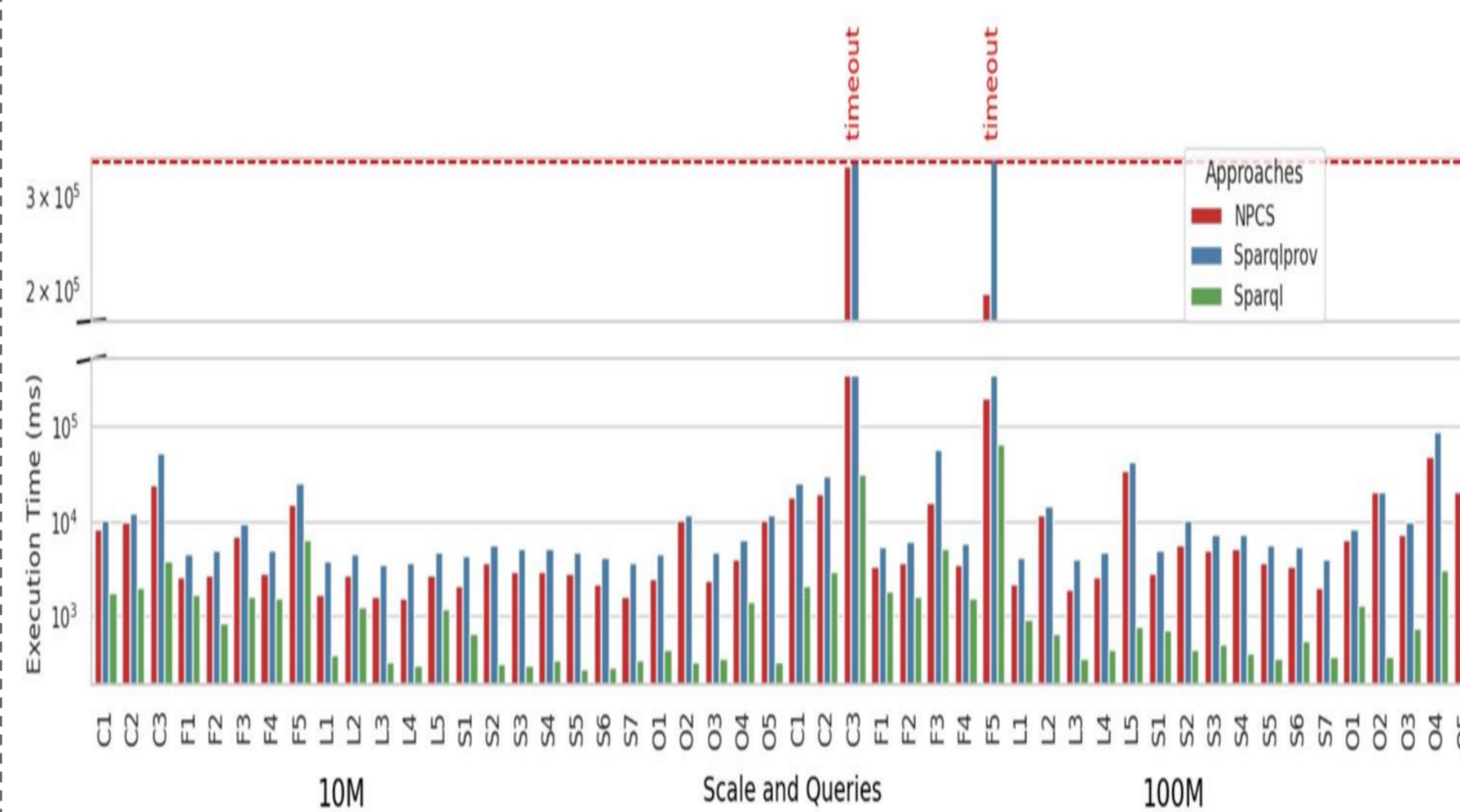
## NPCS Result

[ ?x || k  
 G. Mistral ||  $k_1 \otimes k_3$   
 O. Tokarczuk ||  $k_4 \otimes k_7$  ]

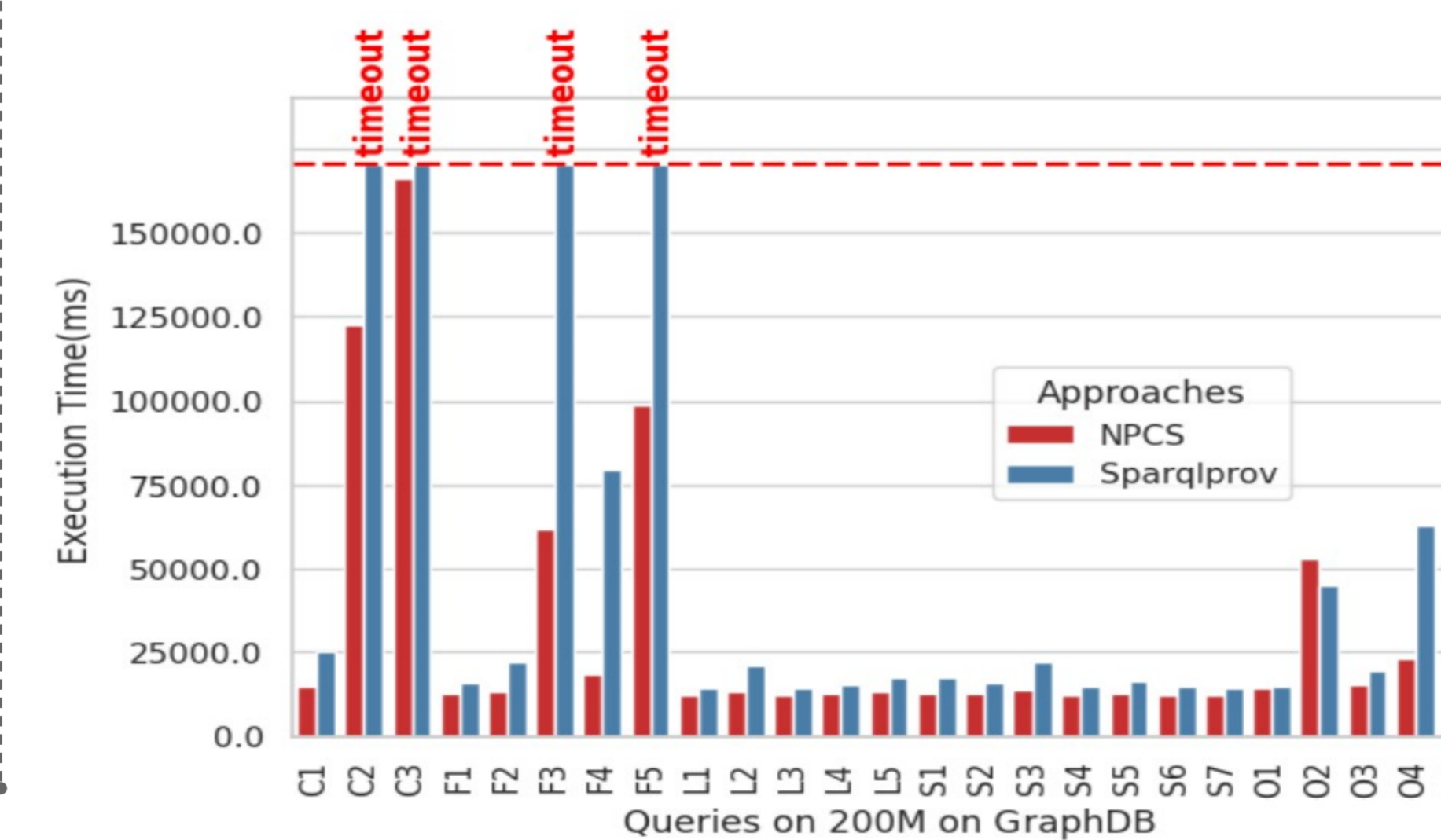
## Evaluation of NPCS

- ❖ Tested on different reification schemes using two engines (Stardog, GraphDB) with datasets of 10M, 100M, 200M from WatDiv and 15 billion triples from Wikidata
- ❖ Observed trends: NPCS consistently outperforms SPARQLprov in 48 out of 50 studied cases
- ❖ Applicable to any standard RDF/SPARQL engine with a significant performance margin

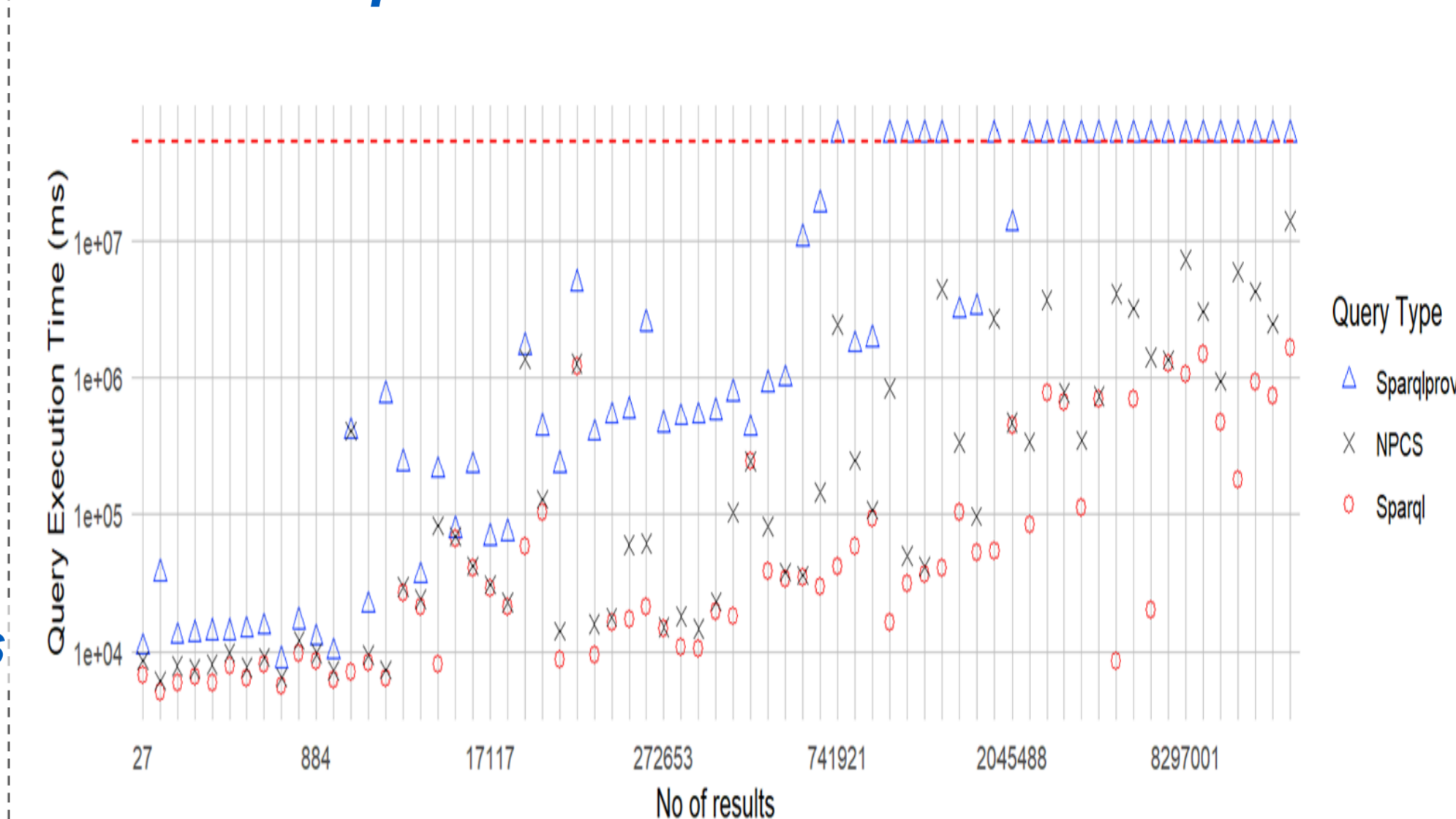
## NPCS performs better on 10M and 100M datasets



## On 200M, NPCS exhibits a more significant performance, with SPARQLprov experiencing increased timeouts



## Graph showing results vs. execution time for WDBench queries on Wikidata



## Conclusion

- ❖ Novel SPARQL-based method for computing how-provenance annotations
- ❖ Outperforms existing solutions
- ❖ Enables efficient computation of provenance for millions of query results on large knowledge graphs
- ❖ Ideal for ETL processes in multi-source KG construction

## Acknowledgment

This work was supported by the ITN KnowGraphs project, under the EU H2020 Marie Skłodowska-Curie grant agreement No 860801

## References

1. Daniel Hernández, Luis Galárraga, and Katja Hose. 2021. Computing How-Provenance for SPARQL Queries via Query Rewriting. Proceedings of the VLDB Endowment 14, 13 (2021), 3389–340.
2. Floris Geerts, Thomas Unger, Grigoris Karvounarakis, Irini Fundulaki, and Vassilis Christophides. 2016. Algebraic Structures for Capturing the Provenance of SPARQL Queries. Journal of the ACM 63, 1 (2016), 7:1–7:63.