

# CHIST-ERA Triple: improving data interoperability and federation across RDF knowledge graphs and Solid Pods

Jerven Bolleman<sup>1</sup>, Elias Crum<sup>2</sup>, Iulian Dragan<sup>1</sup>, Jakub Galgonek<sup>3</sup>, Mark Ibberson<sup>1</sup>, Tarcisio Mendes de Farias<sup>1</sup>, Marek Moos<sup>3</sup>, Marco Pagni<sup>1</sup>, Ruben Taelman<sup>2</sup>, Jiří Vondrášek<sup>3</sup> and Ana Claudia Sima<sup>1,\*</sup>

<sup>1</sup>SIB Swiss Institute of Bioinformatics

<sup>2</sup>Ghent University

<sup>3</sup>IOCB Institute of Organic Chemistry and Biochemistry of the Czech Academy of Sciences

## Abstract

The TRIPLE project, a collaborative effort between the SIB Swiss Institute of Bioinformatics, the University of Ghent and the IOCB Prague, aims to boost the (re)usability of existing knowledge graph resources and improve software tools for RDF data access, documentation and data model visualization. In addition, TRIPLE will increase interoperability between existing public SPARQL endpoints and private data stored in Solid Pods, thus creating an ecosystem of research data that can be seamlessly integrated through efficient and expressive federated SPARQL queries.

## Keywords

RDF, federated SPARQL, Solid Pods, Open Research Data

The Resource Description Framework (RDF) provides a powerful way to structure data resources, coupled with the SPARQL query language, which can be used to interrogate these data, even when they are physically distributed, through the use of federated queries. Although powerful, these technologies are still currently underexploited due to significant technical and usability barriers limiting their use to a select few researchers. Knowledge graphs in RDF are generally not sufficiently described, with limited documentation, making it difficult to ascertain what to query and how to do so. In addition, execution times for complex queries are often very long and difficult to optimise. Finally, the query results are often difficult to integrate with

---

SWAT4HCLS 2024: Bridging Life Sciences and Technology, February 26-29, Leiden, The Netherlands


\*Corresponding author.

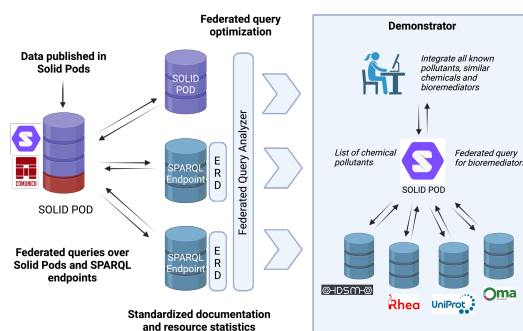
✉ jerven.bolleman@sib.swiss (J. Bolleman); elias.crum@ugent.be (E. Crum); iulian.dragan@sib.swiss (I. Dragan); jakub.galgonek@uochb.cas.cz (J. Galgonek); mark.ibberson@sin.swiss (M. Ibberson); tarcisio.medes@sib.swiss (T. M. d. Farias); marek.moos@uochb.cas.cz (M. Moos); marco.pagni@sib.swiss (M. Pagni); ruben.taelman@UGent.be (R. Taelman); jiri.vondrasek@uochb.cas.cz (J. Vondrášek); ana-claudia.sima@sib.swiss (A. C. Sima)

ORCID 0000-0002-7449-1266 (J. Bolleman); 0009-0005-3991-754X (E. Crum); 0000-0002-7038-544X (J. Galgonek); 0000-0003-3152-5670 (M. Ibberson); 0000-0002-3175-5372 (T. M. d. Farias); 0009-0008-9770-3971 (M. Moos); 0000-0001-9292-9463 (M. Pagni); 0000-0001-5118-256X (R. Taelman); 0000-0002-6066-973X (J. Vondrášek); 0000-0003-3213-4495 (A. C. Sima)



© 2024 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

 CEUR Workshop Proceedings (CEUR-WS.org)



**Figure 1:** TRIPLE system architecture and data flow, interoperating Solid Pods and SPARQL endpoints

private data, such as preliminary, unpublished results. The TRIPLE project, a collaborative effort between the SIB Swiss Institute of Bioinformatics, the University of Ghent and the IOCB Prague, will address these challenges by developing innovative solutions on four fronts:

1. Storing private (unpublished) data in Solid [1] Pods, an emerging technology enabling decentralised data vaults to host private RDF endpoints, execute federated SPARQL queries and cache results.
2. Optimising federated queries spanning public and private SPARQL endpoints allowing users to query multiple resources from within their Solid Pod.
3. Adapting state-of-the-art RDF documentation tools and making them available for all SPARQL endpoints, including Solid Pods.
4. Developing data model visualisations, sets of standardised federated queries, and query analysis tools to help users understand the data and run efficient federated queries.

Finally, a demonstrator will show the impact of these advances when applied to a technically challenging use case of scientific relevance: the search for suitable organisms for bioremediation<sup>1</sup>. The components and data flow in the TRIPLE system architecture are illustrated in Figure 1. The TRIPLE project will improve the (re)usability of existing knowledge graph resources and software tools to access them. In doing so, TRIPLE will create conditions for reproducible research in any domain based on open or shared data and software. Furthermore, the interoperability with Solid Pods will enable researchers to integrate their data with knowledge graphs.

## Acknowledgments

We acknowledge support from the CHIST-ERA Open Research Data (ORD) grant. The SIB received funding from the Swiss National Science Foundation (SNSF). The University of Ghent acknowledges funding from the Research Foundation – Flanders (FWO). The IOCB Prague is grateful for funding from the Technology Agency of the Czech Republic (TAČR) within the National Recovery Plan, project No. TH86010003.

<sup>1</sup><https://en.wikipedia.org/wiki/Bioremediation>

## References

- [1] A. V. Sambra, E. Mansour, S. Hawke, M. Zereba, N. Greco, A. Ghanem, D. Zagidulin, A. Abounaga, T. Berners-Lee, Solid: a platform for decentralized social applications based on linked data, MIT CSAIL & Qatar Computing Research Institute, Tech. Rep. (2016).