

# Unveiling Scholarly Communities over Knowledge Graphs

Sahar Vahdati<sup>1</sup>[0000-0002-7171-169X], Guillermo Palma<sup>2</sup>[0000-0002-8111-2439],  
Rahul Jyoti Nath<sup>1</sup>, Christoph Lange<sup>1,4</sup>[0000-0001-9879-3827], Sören  
Auer<sup>2,3</sup>[0000-0002-0698-2864], and Maria-Esther Vidal<sup>2,3</sup>[0000-0003-1160-8727]

<sup>1</sup> University of Bonn, Germany

{vahdati, lange}@cs.uni-bonn.de, s6ranath@uni-bonn.de,

<sup>2</sup> L3S Research Center, Germany {palma, auer, vidal}@L3S.de

<sup>3</sup> TIB Leibniz Information Centre for Science and Technology, Hannover, Germany  
Maria.Vidal@tib.eu

<sup>4</sup> Fraunhofer IAIS, Germany

**Abstract.** Knowledge graphs represent the meaning of properties of real-world entities and relationships among them in a natural way. Exploiting semantics encoded in knowledge graphs enables the implementation of knowledge-driven tasks such as semantic retrieval, query processing, and question answering, as well as solutions to knowledge discovery tasks including pattern discovery and link prediction. In this paper, we tackle the problem of knowledge discovery in scholarly knowledge graphs, i.e., graphs that integrate scholarly data, and present KORONA, a knowledge-driven framework able to unveil scholarly communities for the prediction of scholarly networks. KORONA implements a graph partition approach and relies on semantic similarity measures to determine relatedness between scholarly entities. As a proof of concept, we built a scholarly knowledge graph with data from researchers, conferences, and papers of the Semantic Web area, and apply KORONA to uncover co-authorship networks. Results observed from our empirical evaluation suggest that exploiting semantics in scholarly knowledge graphs enables the identification of previously unknown relations between researchers. By extending the ontology, these observations can be generalized to other scholarly entities, e.g., articles or institutions, for the prediction of other scholarly patterns, e.g., co-citations or academic collaboration.

## 1 Introduction

Knowledge semantically represented in knowledge graphs can be exploited to solve a broad range of problems in the respective domain. For example, in scientific domains, such as bio-medicine, scholarly communication, or even in industries, knowledge graphs enable not only the description of the meaning of data, but the integration of data from heterogeneous sources and the discovery of previously unknown patterns. With the rapid growth in the number of publications, scientific groups, and research topics, the availability of scholarly datasets has considerably increased. This generates a great challenge for researchers, particularly, to keep track of new published scientific results and potential future co-authors. To alleviate the impact of the explosion of scholarly

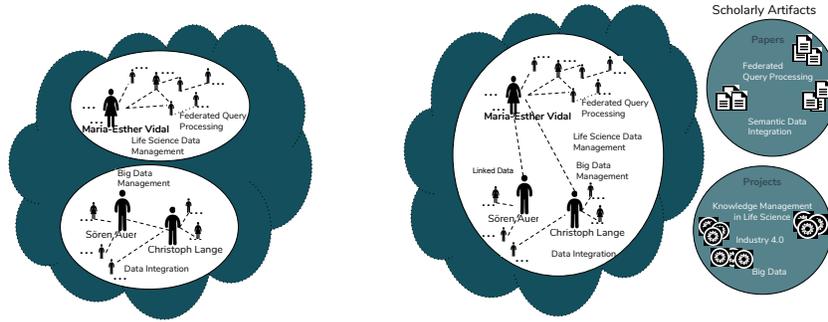
data, knowledge graphs provide a formal framework where scholarly datasets can be integrated and diverse knowledge-driven tasks can be addressed. Nevertheless, to exploit the semantics encoded in such knowledge graphs, a deep analysis of the graph structure as well as the semantics of the represented relations, is required. There have been several attempts considering both of these aspects. However, the majority of previous approaches rely on the topology of the graphs and usually omit the encoded meaning of the data. Most of such approaches are also mainly applied on special graph topologies, e.g., ego networks rather than general knowledge graphs. To provide an effective solution to the problem of representing scholarly data in knowledge graphs, and exploiting them to effectively support knowledge-driven tasks such as pattern discovery, we propose KORONA, a knowledge-driven framework for scholarly knowledge graphs. KORONA enables both the creation of scholarly knowledge graphs and knowledge discovery. Specifically, KORONA resorts to community detection methods and semantic similarity measures to discover hidden relations in scholarly knowledge graphs. We have empirically evaluated the performance of KORONA in a knowledge graph of publications and researchers from the Semantic Web area. As a proof of concept, we studied the accuracy of identifying co-author networks. Further, the predictive capacity of KORONA has been analyzed by members of the Semantic Web area. Experimental outcomes suggest the next conclusions: *i*) KORONA identifies co-author networks that include researchers that both work on similar topics, and attend and publish in the same scientific venues. *ii*) KORONA allows for uncovering scientific relations among researchers of the Semantic Web area. The contributions of this paper are as follows:

- A scholarly knowledge graph integrating data from DBLP datasets;
- The KORONA knowledge-driven framework, which has been implemented on top of two graph partitioning tools, semEP [8] and METIS [3], and relies on semantic similarity to identify patterns in a scholarly knowledge graph;
- Collaboration suggestions based on co-author networks; and
- An empirical evaluation of the quality of KORONA using semEP and METIS.

This paper includes five additional sections. Section 2 motivates our work with an example. The KORONA approach is presented in section 3. Related work is analyzed in section 4. Section 5 reports on experimental results. Finally, section 6 concludes and presents ideas for future work.

## 2 Motivating Example

In this section, we motivate the problem of knowledge discovery tackled in this paper. We present an example of co-authorship relation discovery between researchers working on data-centric problems in the Semantic Web area. We checked the Google Scholar profiles of three researchers between 2015 and 2017, and compared their networks of co-authorship. By 2016, Sören Auer and Christoph Lange were part of the same research group and wrote a large number of joint publications. Similarly, Maria-Esther Vidal, also working on data management



(a) Researchers working on similar topics were in two co-authorship communities.

(b) Researchers working on similar topics constitute a co-authorship community and produce a large number of scholarly artifacts.

Fig. 1: **Motivating Example.** Co-authorship communities from the Semantic Web area working on data-centric problems. Researchers were in different co-authorship communities (2016) (a) started a successful scientific collaboration in 2016 (b), and as a result, produced a large number of scholarly artifacts.

topics, was part of a co-authorship community. Figure 1b illustrates the two co-authorship communities, which were confirmed by the three researchers. After 2016, these three researchers started to work in the same research lab, and a large number of scientific results, e.g., papers and projects, was produced. An approach able to discover such potential collaborations automatically would allow for the identification of the best collaborators and, thus, for maximizing the success chances of scholars and researchers working on similar scientific problems. In this paper, we rely on the natural intuition that successful researchers working on similar problems and producing similar solutions can collaborate successfully, and propose KORONA, a framework able to discover unknown relations between scholarly entities in a knowledge graph. KORONA implements graph partitioning methods able to exploit semantics encoded in a scholarly knowledge graph and to identify communities of scholarly entities that should be connected or related.

### 3 Our Approach: Korona

#### 3.1 Preliminaries

The definitions required to understand our approach are presented in this section. First, we define a scholarly knowledge graph as a knowledge graph where nodes represent scholarly entities of different types, e.g., publications, researchers, publication venues, or scientific institutions, and edges correspond to an association between these entities, e.g., co-authors or citations.

**Definition 1** *Scholarly Knowledge Graph.* Let  $U$  be a set of RDF URI references and  $L$  a set of RDF literals. Given sets  $V_e$  and  $V_t$  of scholarly entities and types, respectively, and given a set  $P$  of properties representing scholarly relations, a scholarly knowledge graph is defined as  $SKG = (V_e \cup V_t, E, P)$ , where:

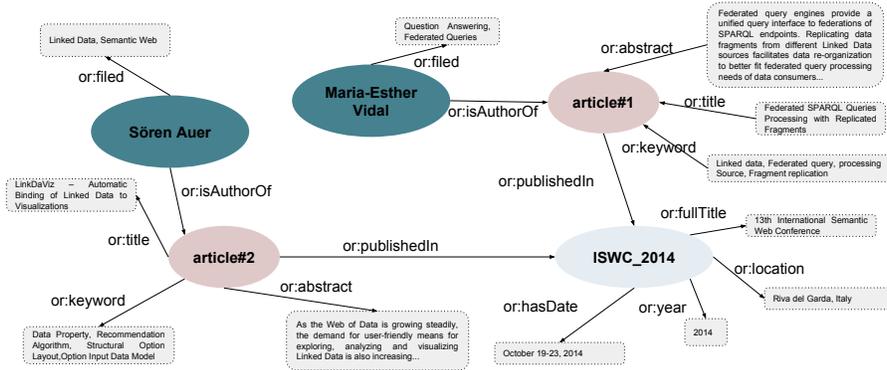


Fig. 2: **Korona Knowledge Graph.** Scholarly entities and relations.

- Scholarly entities and types are represented as RDF URIs, i.e.,  $V_e \cup V_t \subseteq U$ ;
- Relations between scholarly entities and types are represented as RDF properties, i.e.,  $P \subseteq U$  and  $E \subseteq (V_e \cup V_t \times P \times V_e \cup V_t \cup U)$

Figure 2 shows a portion of a scholarly knowledge graph describing scholarly entities, e.g., papers, publication venues, researchers, and different relations among them, e.g., co-authorship, citation, and collaboration.

**Definition 2** *Co-author Network.* A co-author network  $\mathcal{CAN}=(V_a, E_a, P_a)$  corresponds to a subgraph of  $SKG=(V_e \cup V_t, E, P)$ , where

- Nodes are scholarly entities of type researcher,

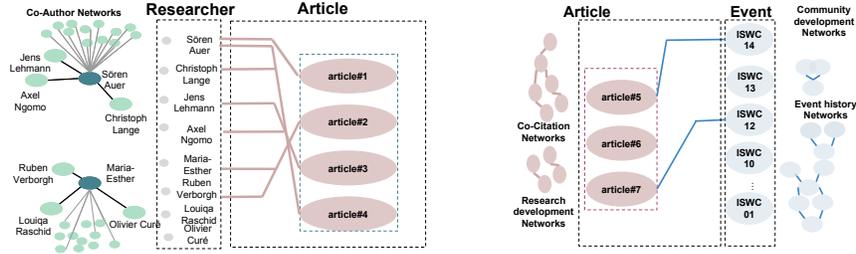
$$V_a = \{a \mid (a \text{ rdf:type :Researcher}) \in E\}$$

- Researchers are related according to co-authorship of scientific publications,  $E_a = \{(a_i : \text{co-author } a_j) \mid \exists p . a_i, a_j \in V_a \wedge (a_i : \text{author } p) \in E \wedge (a_j : \text{author } p) \in E \wedge (p \text{ rdf:type :Publication}) \in E\}$

Figure 3 shows scholarly networks that can be generated by KORONA. Some of these networks are among the recommended applications for scholarly data analytics in [14]. However, the focus on this work is on co-author networks.

### 3.2 Problem Statement

Let  $SKG'=(V_e \cup V_t, E', P)$  and  $SKG=(V_e \cup V_t, E, P)$  be two scholarly knowledge graphs, such that  $SKG'$  is an *ideal* scholarly knowledge graph that contains all the *existing and successful relations* between scholarly entities in  $V_e$ , i.e., an oracle that knows whether two scholarly entities should be related or not.  $SKG=(V_e \cup V_t, E, P)$  is the *actual* scholarly knowledge graph, which only contains a portion of the relations represented in  $SKG'$ , i.e.,  $E \subseteq E'$ ; it represents those relations that are known and is not necessarily complete. Let  $\Delta(E', E) = E' - E$  be the set of relations existing in the ideal scholarly knowledge graph  $SKG'$  that are not represented in the actual scholarly knowledge graph  $SKG$ . Let



(a) Network of Researchers and Articles. (b) Networks of Events and Articles.

Fig. 3: **Scholarly networks.** (a) Co-authors networks from researchers and articles. (b) Co-citation networks from discovered from events and articles.

$SKG_{\text{comp}}=(V_e \cup V_t, E_{\text{comp}}, P)$  be a *complete* knowledge graph, which includes a relation for each possible combination of scholarly entities in  $V_e$  and properties in  $P$ , i.e.,  $E \subseteq E' \subseteq E_{\text{comp}}$ . Given a relation  $e \in \Delta(E_{\text{comp}}, E)$ , the problem of discovering scholarly relations consists in determining whether  $e \in E'$ , i.e., whether a relation  $r=(e_i p e_j)$  corresponds to an existing relation in the ideal scholarly knowledge graph  $SKG'$ .

In this paper, we specifically focus on the problem of discovering *successful co-authorship relations* between researchers in scholarly knowledge graph  $SKG=(V_e \cup V_t, E, P)$ . Thus, we are interested in finding the co-author network  $CAN=(V_a, E_a, P_a)$  composed of the maximal set of relationships or edges that belong to the ideal scholarly knowledge graph, i.e., the set  $E_a$  in  $CAN$  that corresponds to a solution of the following optimization problem:

$$\operatorname{argmax}_{E_a \subseteq E_{\text{comp}}} |E_a \cap E'| \quad (1)$$

### 3.3 Proposed Solution

We propose KORONA to solve the problem of discovering meaningful co-authorship relations between researchers in scholarly knowledge graphs. KORONA relies on information about relatedness between researchers to identify communities composed of researchers that work on similar problems and publish in similar scientific events. KORONA is implemented as an unsupervised machine learning method able to partition a scholarly knowledge graph into subgraphs or communities of co-author networks. Moreover, KORONA applies the *homophily* prediction principle over the communities of co-author networks to identify successful co-author relations between researchers in the knowledge graph. The *homophily* prediction principle states that similar entities tend to be related to similar entities [6]. Intuitively, the application of the *homophily* prediction principle enables KORONA to relate two researchers  $r_i$  and  $r_j$  whenever they work on similar research topics or publish in similar scientific venues. The relatedness or similarity between two scholarly entities, e.g., researchers, research topics, or scientific venues, is represented as RDF properties in the scholarly knowledge graph. Semantic similarity measures, e.g., GADES [10] or Doc2Vec [5], are utilized to

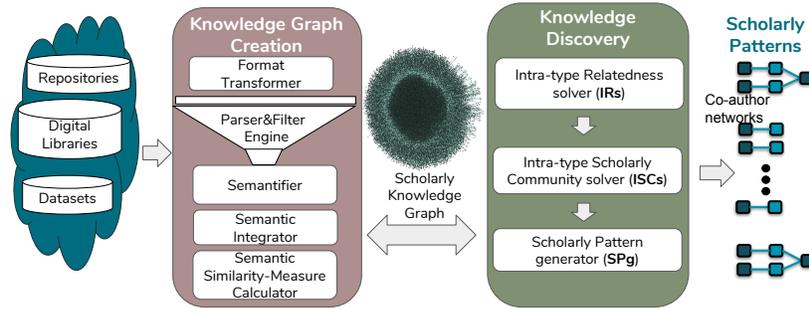


Fig. 4: **The KORONA Architecture.** KORONA receives scholarly datasets and outputs scholarly patterns, e.g., co-author networks. First, a scholarly knowledge graph is created. Then, community detection methods and similarity measures are used to compute communities of scholarly entities and scholarly patterns.

quantify the degree of relatedness between two scholarly entities. The identified degree shows the relevance of entities and returns the most related ones.

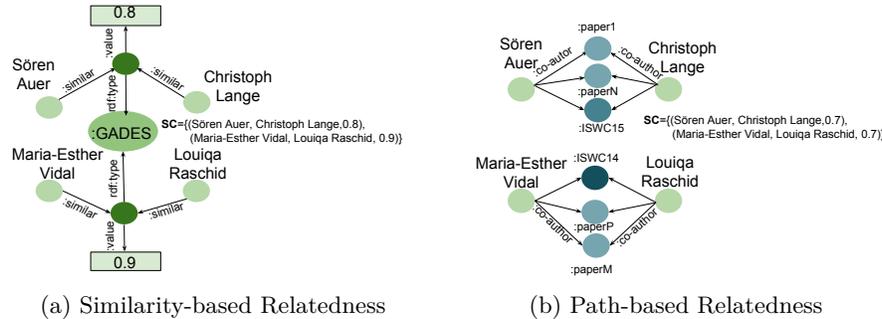


Fig. 5: **Intra-type Relatedness solver (IRs).** Relatedness across scholarly entities. (a) Relatedness is computed according to the values of a semantic similarity metrics, e.g., GADES. (b) Relatedness is determined based on the number of paths between two scholarly entities.

Figure 4 depicts the KORONA architecture; it implements a knowledge-driven approach able to transform scholarly data ingested from publicly available data sources into patterns that represent discovered relationships between researchers. Thus, KORONA receives scholarly data sources and outputs co-author networks; it works in two stages: (a) Knowledge graph creation and (b) Knowledge graph discovery. During the knowledge graph creation stage, a semantic integration pipeline is followed in order to create a scholarly knowledge graph from data ingested from heterogeneous scholarly data sources. It utilizes mapping rules between the KORONA ontology and the input data sources to create the scholarly knowledge graph. Additionally, semantic similarity measures are used to compute

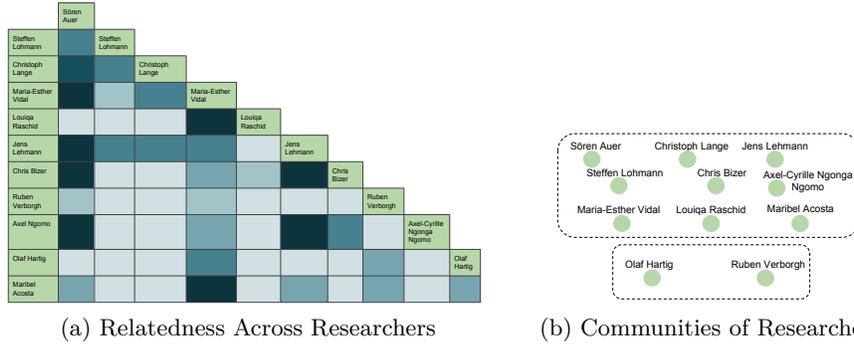


Fig. 6: **Intra-type Relatedness solver (IRs)**. Communities of similar researchers are computed. (a) The tabular representation of  $\mathcal{SC}$ ; lower and higher values of similarity are represented by lighter and darker colors, respectively. (b) Two communities of researchers; each one includes highly similar researchers.

the relatedness between scholarly entities; the results are explicitly represented in the knowledge graph as scores in the range of 0.0 and 1.0. The knowledge graph creation stage is executed offline and enables the integration of new entities in the knowledge graph whenever the input data sources change. On the other hand, the knowledge graph discovery step is executed *on the fly* over an existing scholarly knowledge graph. During this stage, KORONA executes three main tasks: (i) Intra-type Relatedness solver (**IRs**); (ii) Intra-type Scholarly Community solver (**IRSCs**); and (iii) Scholarly Pattern generator (**SPg**).

**Intra-type Relatedness solver (IRs)**. This module quantifies relatedness between the scholarly entities of the same type in a scholarly knowledge graph  $SKG=(V_e \cup V_t, E, P)$ . **IRs** receives as input  $SKG=(V_e \cup V_t, E, P)$  and a scholarly type  $V_a$  in  $V_t$ ; it outputs a set  $\mathcal{SC}$  of triples  $(e_i, e_j, score)$ , where  $e_i$  and  $e_j$  belong to  $V_a$  and  $score$  quantifies the relatedness between  $e_i$  and  $e_j$ . The relatedness can be just computed in terms of the values of similarity represented in the knowledge graph, e.g., according to the values of the semantic similarity according to GADES or Doc2Vec. Alternatively, the values of relatedness can be computed based on the number of paths in the scholarly knowledge graph that connect the scholarly entities  $e_i$  and  $e_j$ . Figure 5 depicts two representations of the relatedness of scholarly entities. As shown in Figure 5a, **IRs** generates a set  $\mathcal{SC}$  according to the GADES values of semantic similarity; thus, **IRs** includes two triples (Sören Auer, Christoph Lange, 0.8), (Maria-Esther Vidal, Louïqa Raschid, 0.9) in  $\mathcal{SC}$ . On the other hand, if paths between scholarly entities are considered (Figure 5b), the values of relatedness can differ, e.g., in this case, Sören Auer and Christoph Lange are equally similar as Maria-Esther Vidal and Louïqa Raschid.

**Intra-type Scholarly Community solver (IRSCs)**. Once the relatedness between the scholarly entities has been computed, communities of highly related scholarly entities are determined. **IRSCs** resorts to unsupervised methods such as METIS or semEP, and to relatedness values stored in  $\mathcal{SC}$ , to compute the scholarly communities. Figure 6 depicts scholarly communities computed by **IRSCs** based on similarity values; as observed, each community includes

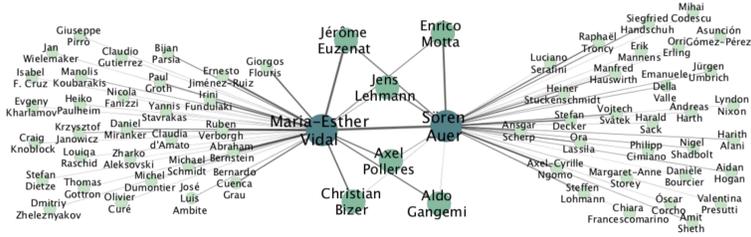


Fig. 7: **Co-author network.** A network generated from scholarly communities.

researchers that are highly related; for readability,  $\mathcal{SC}$  is shown as a heatmap where lower and higher values of similarity are represented by lighter and darker colors, respectively. For example, in Figure 6a, Sören Auer, Christoph Lange, and Maria-Esther Vidal are quite similar, and they are in the same community.

**Scholarly Pattern generator (SPg).** **SPg** receives communities of scholarly entities and produces a network, e.g., a co-author network. **SPg** applies the *homophily* prediction principle on the input communities, and connects the scholarly entities in one community in a network. Figure 7 shows a co-author network computed based on a scholarly knowledge graph created from DBLP; as observed, Sören Auer, Christoph Lange, and Maria-Esther Vidal are included in the same co-author network. In addition to computing the scholarly networks, **SPg** scores the relations in a network and computes the *weight of connectivity* of a relation between two entities. For example, in Figure 7, thicker lines represent strongly connected researchers in the network. **SPg** can also filter from a network the relations labeled with higher values of weight of connectivity. All the relations in a network correspond to solutions to the problem of discovering *successful co-authorship relations* defined in Equation 1. To compute the weights of connectivity, **SPg** considers the values of similarity of the scholarly entities in a community  $C$ ; weights are computed as aggregated values using an aggregation function  $f(\cdot)$ , e.g., average or triangular norm. For each pair  $(e_i, e_j)$  of scholarly entities in  $C$ , the weight of connectivity between  $e_i$  and  $e_j$ ,  $\phi(e_i, e_j | C)$ , is defined as:  $\phi(e_i, e_j | C) = \{f(\text{score}) \mid e_z, e_q \in C \wedge (e_z, e_q, \text{score}) \in \mathcal{SC}\}$ .

## 4 Empirical Evaluation

### 4.1 Knowledge Graph Creation

A scholarly knowledge graph has been crafted using the DBLP collection (7.83 GB in April 2017<sup>5</sup>); it includes researchers, papers, and publication year from the International Semantic Web Conference (ISWC) 2001–2016. The knowledge graph also includes similarity values between researchers who have published at ISWC (2001–2017). Let  $PC_{e_i}$  and  $PC_{e_j}$  be the number of papers published by researchers  $e_i$  and  $e_j$  together (as co-authors), respectively at ISWC (2001–2017). Let  $TP_{e_i}$  and  $TP_{e_j}$  be the total number of papers that  $e_i$  and  $e_j$  have in all conferences of the scholarly knowledge graph, respectively. The similarity measure is defined as:  $\text{Sim}R(e_i, e_j) = \frac{PC_{e_i} \cap PC_{e_j}}{TP_{e_i} \cup TP_{e_j}}$ . The similarities between ISWC

<sup>5</sup> <http://dblp2.uni-trier.de/e55477e3eda3bfd402faefd37c7a8d62/>

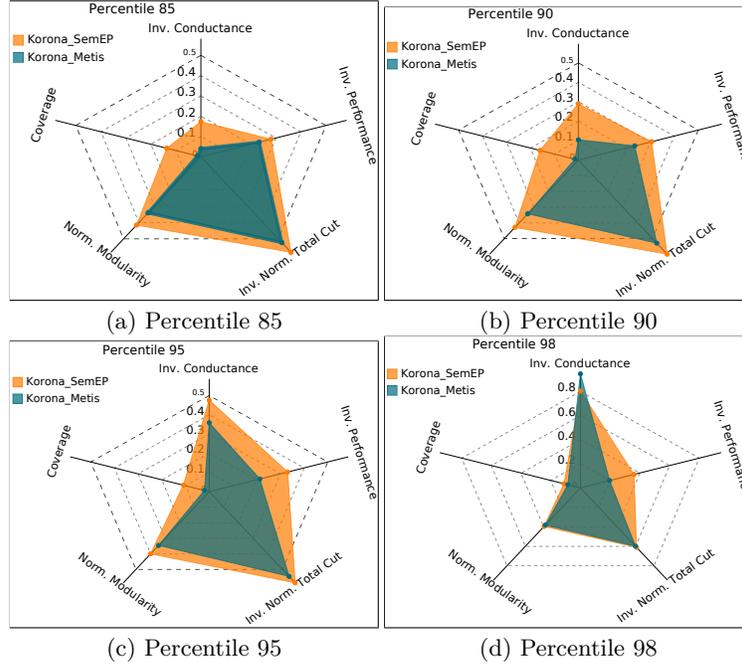


Fig. 8: **Quality of Korona.** Communities evaluated in terms of prediction metrics (higher values are better); percentiles 85, 90, 95, and 98 are reported. KORONA exhibits the best performance at percentile 95 and groups similar researchers according to research topics and events where they publish.

(2002–2016) are represented as well. Let  $RC_i$  and  $RC_j$  the number of the authors with papers published in conferences  $c_i$  and  $c_j$  respectively. The similarity measure corresponds to  $SimC(c_i, c_j) = \frac{RC_i \cap RC_j}{RC_i \cup RC_j}$ . Thus, the scholarly knowledge graph includes both scholarly entities enriched with their values of similarity.

## 4.2 Experimental Study

The effectiveness of KORONA has been evaluated in terms of the quality of both the generated communities of researchers and the predicted co-author networks.

**Research Questions:** We assess the following research questions: **RQ1)** Does the semantics encoded in scholarly knowledge graphs impact the quality of scholarly patterns? **RQ2)** Does the semantics encoded in scholarly knowledge graph allow for improving the quality of the predicting co-author relations?

**Implementation:** KORONA is implemented in Python 2.7. The experiments were executed on a macOS High Sierra 10.13 (64 bits) Apple MacBook Air machine with an Intel Core i5 1.6 GHz CPU and 8 GB RAM. METIS 5.1<sup>6</sup> and SemEP<sup>7</sup> are part of KORONA and used to obtain the scholarly patterns.

<sup>6</sup> <http://glaros.dtc.umn.edu/gkhome/metis/metis/download>

<sup>7</sup> <https://github.com/gpalma/semEP>

---

**Q1. Do you know this person? Have you co-authored before?** To avoid confusion, the meaning of “knowing” was kept simple and general. The participants were asked to only consider if they were aware of the existence of the recommended person in their research community.

---

**Q2. Have you co-authored “before” with this person at any event of the ISWC series?** With the same intent of keeping the survey simple, all types of collaboration on papers in any edition of this event series were considered as “having co-authored before”.

---

**Q3. Have you co-authored with this person after May 2016?** Our study considered scholarly metadata of publications until May 2016. The objective of this question was to find out whether a prediction had actually come true, and the researchers had collaborated.

---

**Q4. Have you ever planned to write a paper with the recommended person and you never made it and why?** The aim is to know whether two researchers who had been predicted to work together actually wanted to but then did not and the reason, e.g., geographical distance.

---

**Q5. On a scale from 1–5, (5 being most likely), how do you score the relevance of your research with this person?** The aim is to discover how close and relevant are the collaboration recommendations to the survey participant.

---

Table 1: **Survey.** Questions to validate the recommended collaborations.

**Evaluation metrics:** Let  $Q = \{C_1, \dots, C_n\}$  be the set of communities obtained by KORONA: *Conductance*: measures relatedness of entities in a community, and how different they are to entities outside the community [2]. The inverse of the conductance  $1 - \text{Conductance}(S)$  is reported. *Coverage*: compares the fraction of intra-community similarities among entities to the sum of all similarities among entities [2]. *Modularity*: is the value of the intra-community similarities among the entities divided by the sum of all the similarities among the entities, minus the sum of the similarities among the entities in different communities, in the case they were randomly distributed in the communities [7]. The value of the modularity lies in the range  $[-0.5, 1]$ , which can be scaled to  $[0, 1]$  by computing  $\frac{\text{Modularity}(Q)+0.5}{1.5}$ . *Performance*: sums the number of intra-community relationships, plus the number of non-existent relationships between communities [2]. *Total Cut*: sums all similarities among entities in different communities [1]. Values of total cut are normalized by dividing by the sum of the similarities among the entities; inverse values are reported, i.e.,  $1 - \text{NormTotalCut}(Q)$ .

### Experiment 1: Evaluation of the Quality of Collaboration Patterns.

Prediction metrics are used to evaluate the quality of the communities generated by KORONA using METIS and semEP; relatedness of the researchers is measured in terms of *SimR* and *SimC*. Communities are built according to different similarity criteria; percentiles of 85, 90, 95, and 98 of the values of similarity are analyzed. For example, in percentile 85 only 85% of all similarity values among entities have scores lower than the similarity value in the percentile 85. Figure 8 presents the results of the studied metrics. In general, in all percentiles, the communities include closely related researchers. However, both implementations of KORONA exhibit quite good performance at percentile 95, and allow for grouping together researchers that are highly related in terms of the research topics on which they work, and the events where their papers are published. On the contrary, KORONA creates many communities of no related authors for percentiles 85 and 90, thus exposing low values of coverage and conductance.

<b>Korona</b>	%	Q.1(a)	Q.1(b)	Q.2	Q.3	Q.4	Q.5
Korona-METIS	85	0.26±0.25	0.72±0.29	0.99±0.04	0.86±0.13	0.86±0.20	3.10±0.59
Korona-semEP	85	0.24±0.21	0.80±0.34	1.00±0.00	0.97±0.07	0.93±0.16	3.35±0.85
Korona-METIS	90	0.39±0.24	0.91±0.19	1.00±0.00	1.00±0.00	0.98±0.04	3.03±0.79
Korona-semEP	90	0.13±0.18	0.89±0.18	1.00±0.00	1.00±0.00	0.85±0.23	3.12±1.06
Korona-METIS	95	0.40±0.34	0.93±0.08	1.00±0.00	0.80±0.45	0.95±0.10	3.20±0.81
Korona-semEP	95	0.14±0.30	0.81±0.40	0.67±0.58	0.60±0.55	0.69±0.47	3.83±0.76

Table 2: **Survey results.** Aggregated normalized values of negative answers provided by the study participants during the validation of the recommended collaborations (Q.1(a), Q.1(b), Q.2, Q.3, and Q.4); average (lower is better) and standard deviation (lower is better) are reported. For Q.5, average and standard deviation of the scale from 1–5 are presented; higher average values are better.

**Experiment 2: Survey of the Quality of the Prediction of Collaborations among Researchers.** Results of an online survey<sup>8</sup> among 10 researchers are reported; half of the researchers are from the same research area, while the other half was chosen randomly. Knowledge subgraphs of each of the participants are part of the KORONA research knowledge graph; predictions are computed from these subgraphs. The predictions for each were laid out in an online spreadsheet along with 5 questions and a comment section. Table 1 lists the five questions that the survey participants were asked to validate the answers, while Table 2 reports on the results of the study. The analysis of results suggests that KORONA predictions represent potentially *successful co-authorship relations*; thus, they provide a solution to the problem tackled in this paper.

## 5 Related Work

Xia et al. [14] provides a comprehensive survey of tools and technologies for scholarly data management, as well as a review of data analysis techniques, e.g., social networks and statistical analysis. However, all the proposals have been made over raw data and knowledge-driven methods were not considered. Wang et al. [13] present a comprehensive survey of link prediction in social networks, while Paulheim [9] presents a survey of methodologies used for knowledge graph refinement; both works show the importance of the problem of knowledge discovery. Traverso-Ribón et al. [12] introduces a relation discovery approach, *KOI*, able to identify hidden links in TED talks; it relies on heterogeneous bipartite graphs and on the link discovery approach proposed in [8]. In this work, Palma et al. present semEP, a semantic-based graph partitioning approach, which was used in the implementation of KORONA-semEP. Graph partitioning of semEP is similar to *KOI* with the difference of only considering isolated entities, whereas *KOI* is desired for ego networks. However, it is only applied to ego networks, whereas KORONA is mainly designed for knowledge graphs. Sachan and Ichise [11] propose a syntactic approach considering dense subgraphs of a co-author network created from the DBLP. They discover relations between authors and propose pairs of researchers belonging to the same community. A link discovery tool is

<sup>8</sup> <https://bit.ly/2ENeg2G>

developed for the biomedical domain by Kastrin et al. [4]. Albeit effective, these approaches focus on the graph structure and ignore the meaning of the data.

## 6 Conclusions and Future Work

KORONA is presented for unveiling unknown relations; it relies on semantic similarity measures to discover hidden relations in scholarly knowledge graphs. Reported and validated experimental results show that KORONA retrieves valuable information that can impact the research direction of a researcher. In the future, we plan to extend KORONA to detect other networks, e.g., affiliation networks, co-citation networks and research development networks. We plan to extend our evaluation over big scholarly datasets and study the scalability of KORONA; further, the impact of several semantic similarity measures will be included in the study. Finally, KORONA will be offered as an online service that will enable researchers to explore and analyze the underlying scholarly knowledge graph.

**Acknowledgement** This work has been partially funded by the EU H2020 programme for the project iASiS (grant agreement No. 727658).

## References

1. Buluç, A., Meyerhenke, H., Safro, I., Sanders, P., Schulz, C.: Recent Advances in Graph Partitioning. Springer, Cham (2016)
2. Gaertler, M.: Clustering. In: Network Analysis: Method. Found.
3. Karypis, G., Kumar, V.: A fast and high quality multilevel scheme for partitioning irregular graphs. *Scientific Computing* (1998)
4. Kastrin, A., Rindfleisch, T.C., Hristovski, D.: Link prediction on the semantic MEDLINE network - an approach to literature-based discovery. In: The Discovery Science Conference (2014)
5. Le, Q.V., Mikolov, T.: Distributed representations of sentences and documents. *CoRR* **abs/1405.4053** (2014)
6. Liben-Nowell, D., Kleinberg, J.: The link-prediction problem for social networks. *JASIST* **58**(7) (2007)
7. Newman, M.E.: Modularity and community structure in networks. *Proceedings of the national academy of sciences* **103**(23) (2006)
8. Palma, G., Vidal, M., Raschid, L.: Drug-target interaction prediction using semantic similarity and edge partitioning. In: ISWC (2014)
9. Paulheim, H.: Knowledge graph refinement: A survey of approaches and evaluation methods. *Semantic Web Journal* **8**(3) (2017)
10. Ribón, I.T., Vidal, M., Kämpgen, B., Sure-Vetter, Y.: GADES: A graph-based semantic similarity measure. In: SEMANTICS (2016)
11. Sachan, M., Ichise, R.: Using semantic information to improve link prediction results in network datasets. *IJET* **2**(4) (2010)
12. Traverso-Ribón, I., Palma, G., Flores, A., Vidal, M.E.: Considering semantics on the discovery of relations in knowledge graphs. In: EKAW (2016)
13. Wang, P., Xu, B., Wu, Y., Zhou, X.: Link prediction in social networks: the state-of-the-art. *Link Prediction in Social Networks(SCIS)* **58**(1) (2015)
14. Xia, F., Wang, W., Bekele, T.M., Liu, H.: Big scholarly data: a survey. *IEEE Big Data* (2017)