

Quality-driven Query Processing over Federated RDF Data Sources

Lars Heling

Institute AIFB, Karlsruhe Institute of Technology (KIT), Germany
heling@kit.edu

Abstract. The integration of data from heterogeneous sources is a common task in various domains to enable data-driven applications. Data sources may range from publicly available sources to sources within data lakes of companies. The added value generated by integrating and analyzing the data greatly depends on the quality of the underlying data. As a result, querying heterogeneous data sources as a way of integrating data enabling such applications needs to consider quality aspects. Quality-driven query processing over RDF data sources aims to study approaches which consider data quality description of the data sources to determine optimal query plans. In contrast to most federated query approaches, in quality-driven query processing the quality of an optimal plan and thus of the retrieved data, not only depends on efficiency typically measured as execution time but also on other quality criteria. In this work, we present the challenges associated with considering multiple quality criteria in federated query processing and derive our problem statement accordingly. We present our research questions to address the problem and the associated hypotheses. Finally, we outline our approach including an evaluation plan and provide preliminary results.

Keywords: Federated Querying, Linked Data, Data Quality, SPARQL

1 Introduction

Driven by Linked Open Data (LOD) initiatives, an increasing amount of data is being published as Linked Data (LD) on the web¹. Due to different origins and publishers of the data, the datasets are heterogeneous with respect to various properties, such as schema, data access, and data quality. Data quality may greatly vary due to the procedure and context in which the data is generated and intended to be used for. For instance, automatically extracted data based on crowd-curated data sources such DBpedia² is likely to yield a different data quality than expert-curated datasets which are used in the life science applications such as Drugbank³ [22]. Consequently, the value generated by applications querying the data from such federations of RDF data sources highly depends

¹<https://lod-cloud.net/>

²<https://wiki.dbpedia.org/>

³<http://drugbank.bio2rdf.org/>

on the quality of the sources. A variety of heterogeneous data sources differing in quality are not merely encountered on the web but are also a phenomenon within business organizations. In the era of Big Data, companies increasingly maintain their data in its native, often schema-less format in large repositories, so-called *data lakes*. As a result, employing data analysis over such data lakes requires querying a plethora of heterogeneous data sources differing in schema and data quality. As the effectiveness of analyses depends on the quality of the underlying data, it is desirable to consider data quality aspects when querying the sources [12]. Hence, providing high quality data is an important task in query processing over heterogeneous data sources and it yields three major questions: *i*) how can the data quality of independent and heterogeneous sources be measured, *ii*) which sources should be included when processing a query to provide answers of high quality, and *iii*) how can the quality of a query plan be assessed to be used as a proxy for the quality of the answers it provides?

Several challenges arise when addressing these questions. First, quality assessment approaches retrieving quality descriptions while considering the dynamics and access restriction (i.e., interface to query the data) of the data sources are required. Moreover, in order to find query plans retrieving high quality answers when processing a query, source selection, query decomposition, and query planning need to consider these quality descriptions as well as the fact that the data sources may contain overlapping, replicated or disjoint data. Finally, in contrast to quality-driven query processing over a single data source, the query planning in a federated scenario needs to consider both quality metrics on the data level (e.g., conciseness) as well as on the data source level (e.g., accessibility).

The proposed doctoral thesis aims to address open research questions related to quality-driven querying processing over federations of RDF data sources. In particular, we address the challenges of *i*) retrieving quality descriptions for data provided via web querying interfaces for RDF, *ii*) determining the quality of a query plan with respect to various quality dimensions, and *iii*) finding a quality maximizing query plan to be executed over the federation.

2 State of the Art

Data quality is essential for a large variety of applications in many domains. Data quality is commonly defined as the “fitness for use” [21] which underpins the importance of context when assessing data quality. Federated RDF query processing investigates approaches for querying a federation of RDF data sources and SPARQL is the de facto language to express such queries [1]. Due to the heterogeneity of data sources involved in such federations, data quality may strongly vary across these sources. Extensions of the SPARQL query language which allow considering additional data properties, such as the trustworthiness of statements, have been proposed and quality-driven query processing has been studied in the area of relational databases. However, thus far few research has focused on quality-driven query processing for federated RDF data sources.

Linked Data Quality. Assessing quality for Linked Data has been addressed by various authors [4,7,22]. In their survey, Zaveri et al. [22] identify 18 data quality dimensions for Linked Data and group them into four major dimensions: accessibility, representation, contextual and intrinsic data quality. For each dimension, there are several data quality metrics which are procedures for measuring a certain dimension. Along those lines, RDF dataset profiling, which includes data quality aspects as well, has been proposed to facilitate dataset discovery and selection for tasks such as distributed querying, search and question answering [4]. However, most current distributed query processing approaches merely consider data summaries which do not include quality information of the data sources to decide on the relevancy of data sources [8,9].

Federated RDF Query Processing. Federated SPARQL query engines provide means for integrating RDF data sources by allowing to query a federation of SPARQL endpoints via a unique interface. The initial version of SPARQL does not provide features for querying federations of data sources, however a myriad of approaches to allow for federated query processing have been proposed [3,6,17,19,18]. The prevalent definition for the result of executing a query over a federation has been that the results should be the same as querying the union of all RDF data provided by the federation members, i.e. *complete* answers [1]. As a result, the proposed approaches aim to provide complete answers while minimizing the processing cost and they devise different strategies to achieve this goal. In contrast, our work aims to address additional quality dimensions besides completeness and especially, the potential trade-off between these dimensions.

Quality-Driven RDF Query Processing. Previous works have focused on enhancing the quality of answers when querying RDF. Darari et al. [5] focus on the quality dimension of data completeness in query processing. They propose the concept of completeness statements which allows to decide whether a query returns complete answers with respect to an ideal graph containing all facts that hold in this world. Acosta et al. [2] propose HARE, a hybrid query processing system leveraging crowd sourcing to improve query results by identifying missing values in a data set and completing them using the crowd. Based on a model for completeness, the system detects parts of the query which potentially yield incomplete results and uses the crowd to augment the missing answers.

SPARQL Extensions. Hartig [10] proposes tSPARQL which is an extension allowing to query trust weighted RDF graphs. Such graphs consist of an RDF graph G and a trust function tv^C mapping all triples in G to a trust value. Furthermore, trust weighted solution mappings assigning a trust value to each solution mapping based on the trust weighted graph. The keywords TRUST AS and ENSURE TRUST are introduced to make use of the assigned trust values within a query. The keywords allow for assigning trust values to variables and for filtering of solution mappings according to their trust value.

Similarly, AnQL [14] is an extension of SPARQL for querying annotated RDFS. Annotated RDFS is an extension of the RDF data model providing means for annotations of triples. An annotated triple is an expression $\pi : \lambda$ with π a triple and λ an annotation value. Hence, it is a more generic approach than the

trust weighted RDF graph by allowing to annotate triples with any meta information, such as trust, temporal validity, or provenance. To allow for querying such annotated graphs, the syntax of SPARQL is extended by annotated triple pattern and the semantics of evaluating an AnQL graph pattern are defined accordingly. Moreover, the keywords ASSIGN, ORDERBY and GROUPBY are introduced for variable assignment, projections, aggregates and solution modifiers.

Quality-driven Query Processing in Databases. Previous work in the area of relational databases has studied how data quality may be considered when querying heterogeneous data sources. Naumann et al. [15] describe a framework that includes quality information in query processing over multiple databases. The goal of their work is using the quality information about the data sources in the query planning process and they argue that “the goodness of a plan depends primarily on the expected quality of the results and not on pure technical criteria such as response time” [15]. First, the authors determine a set of information quality measures which they categorize into *source-specific*, *query correspondence assertion(QCA⁴)-specific* and *user query-specific* criteria. Thereafter, these source-specific quality metrics are used to filter out sources in the source selection phase and the QCAs are used in the plan creation phase to create all plans for a user query. Lastly, in the plan selection phase, each plan is given an overall quality value by first calculating the quality values for the QCA-specific and user query-specific criteria, and by finally aggregating the information quality of a plan into a single value using simple additive weighting.

Summarizing, Linked Data quality assessment and federated query processing have been studied extensively in the Semantic Web community. In addition, enhancing the quality of query answers by identifying and completing missing information has also been studied. However, most current approaches aim to optimize answer completeness but do not provide a generic model for the quality of query plans. In addition, quality assessment has not considered restrictions that are relevant in a federated query processing scenario such as the dynamics of data sources and the restriction to access the data. Moreover, approaches to enable quality-driven query processing in heterogeneous information systems covering various quality dimensions have been proposed for relational databases. Due to the semantics of SPARQL and the specifics of its operators, these approaches cannot be transferred seamlessly to querying federated RDF data sources.

3 Problem Statement and Contributions

In federated query processing over RDF data, query plans are determined in three steps: query decomposition, source selection, and query optimization. A query plan P_Q for a SPARQL query Q is a tree where the leaves are requesting a sub-expression of the query at a member of the federation and the inner nodes are SPARQL algebra operators. Furthermore, we denote the universe of all plans for a query Q by \mathcal{P}_Q .

⁴The mediator uses query correspondence assertions (QCAs) in order to determine contents, i.e. available relations, of the sources.

Given a query, it is a difficult task to devise optimal query plans which yield complete answers while minimizing processing cost. This is due to the complexity of exploring the space of plans which is an NP-complete problem [13] as well as the difficulty of estimating the cost of those plans [16,20]. Enabling quality-driven federated query processing over RDF data sources also entails additional challenges: *i*) gathering quality information from various data quality dimensions from the members of the federation, *ii*) quality estimation of query plans and *iii*) determining the best query plan according to the quality criteria. These observations lead to our problem statement.

Problem Statement: Given a SPARQL query Q , a federation of RDF data sources $F = \{D_1, \dots, D_n\}$ and a set of quality descriptions $\mathcal{Q} = \{Q_1, \dots, Q_n\}$ with $Q_i \in \mathcal{Q}$ a quality description for $D_i \in F$. The quality-driven query plan optimization problem is defined as devising a plan P_Q^* such that

$$P_Q^* = \arg \max_{P_Q \in \mathcal{P}_Q} \text{quality}_{\mathcal{Q}}(P_Q)$$

with $\text{quality}_{\mathcal{Q}} : \mathcal{P}_Q \rightarrow \mathbb{R}^+$ a quality assessment function for a set of quality description \mathcal{Q} associating every query plan with a quality value.

In other words, we want to find a plan for a given query which maximizes quality according to the quality descriptions of the data sources in the federation. Please note that we do not explicitly encode answer completeness as a constraint in our problem definition since part of the work aims to investigate the trade-off between answer completeness and other quality dimensions of the query plans. Based on the problem statement and according to the challenges associated with it, we formulate the following research questions:

RQ1 Data Source Quality:

- What quality dimensions of federated RDF data sources can be effectively measured on a fine-grained level?

RQ2 Query Plan Quality:

- How can we measure the quality of a query plan based on the quality descriptions of the data sources with multiple quality criteria?

RQ3 Query Plan Optimization:

- What is the impact on time complexity in the federated query optimization problem when considering multiple quality criteria?
- Which methods are suitable to determine quality maximizing query plans while simultaneously considering execution cost?

Our hypotheses are derived directly from the research questions.

H1 Gathering data quality information from federated RDF data sources and deriving a quality-enhanced source descriptions can be achieved by applying existing data quality metrics and leveraging web query interfaces to RDF data (i.e. SPARQL endpoints and Triple Pattern Fragments).

- H2** Existing cost estimation models for query planning can be extended to incorporate multiple quality criteria and determine an aggregated quality value for query plans which considers both data source and query-level quality criteria.
- H3** Meta-heuristics can be employed to determine (near-)optimal query plans considering various data quality dimensions and therefore, provide query plans with higher quality than existing solutions.

The research to be conducted and the resulting contributions aim to address the presented research problem by investigating the associated hypotheses.

4 Research Methodology and Approach

The methodology adopted in the development of this doctoral work adheres to the following tasks:

1. Investigation of state of the art research relevant to the identified problem. This includes the study of literature about quality assessment and quality-driven query processing in the areas of Databases and Semantic Web.
2. Formalization of the research problem and formulation of research questions and hypotheses.
3. Definition of solutions to address the hypotheses. Identification of novel contributions entailed by the solutions as well as studying formal properties of the proposed solutions.
4. Implementation of the solutions and empirical evaluation of their performance. The performance evaluation is conducted with respect to the state of the art solutions if available. The experiments will be conducted as follows:
 - i*) Implementation of state of the art or baselines approaches.
 - ii*) Designing an evaluation approach according to the studied research question. Reuse and if necessary adjustment of existing benchmark and evaluation criteria.
 - iii*) Execution of experiments based on the benchmarks to obtain data for drawing conclusions about the hypotheses.
 - iv*) Analysis of the results.

In order to address the presented research problem, our approach aims to study each of the presented hypotheses. Figure 1 provides an overview of the approach with the corresponding research questions indicated. Hypothesis **H1** will be addressed by developing methods to assess various data quality dimensions of RDF data sources, i.e. retrieving \mathcal{Q} for F . In contrast to existing methods which typically examine a dump of the entire data available at a data source, our methods will be restricted to data access of web querying interfaces for RDF, namely SPARQL endpoints and Triple Pattern Fragment (TPF) servers. Therefore, we need to select quality dimensions relevant in the realm of query processing and adapt methods to measure these dimensions while considering the access constraints. For instance, the timeliness of the data may be relevant as sources may provide update their data and sources with outdated data may be pruned in

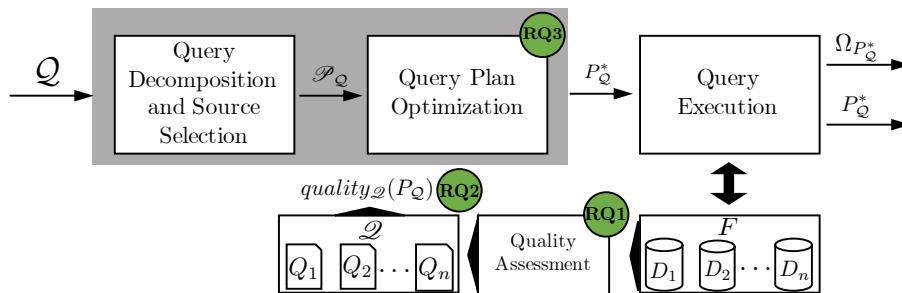


Fig. 1: Our approach: Given a query Q , the quality maximizing query plan $P_Q^* \in \mathcal{P}_Q$ is determined according to the quality description in \mathcal{Q} . The plan is executed over the federation and the results $\Omega_{P_Q^*}$ and optimized plan P_Q^* are returned.

the query plan. Furthermore, we aim to investigate which dimensions can be measured on the data source level and which are query-specific.

We want to address the second hypothesis **H2** by determining approaches to estimate the quality of query plans based on the quality assessment of the previous step, i.e., study how to compute $quality_{\mathcal{Q}}(P_{\mathcal{Q}})$. Existing cost estimation methods for query planning (e.g., [16,20]) should be used as a basis and adapted accordingly. Moreover, we aim to investigate how preferences of a user posing a query may be encoded in quality estimation. In experimental studies, the implementations of the approaches will be used to evaluate their feasibility.

Finally, to test our last hypothesis **H3**, we aim to study different meta-heuristics and adapt them to the problem of determining the quality maximizing query plan P_Q^* . As previously mentioned, the complexity of determining an optimal plan is NP-complete and optimizing for its quality depending on multiple criteria is likely to disallow using existing heuristics. Implementing these heuristics and conducting an experimental study on their performance with respect to relevant metrics, will allow determining the most suitable method.

5 Evaluation Plan

For the evaluation of the presented approach, we aim to conduct both theoretical and experimental evaluations of the proposed solutions. The goal of the theoretical evaluation will be investigating the complexity of query planning when considering several criteria with respect to the quality of the plans. In the experimental evaluation, the hypotheses will be evaluated by implementing the proposed solutions, conducting a series of experiments and analyzing the results.

Query Benchmark and Data Sources: For the experimental evaluation, we require a set of SPARQL queries which can be evaluated over a federation of data sources. Ideally, the data sources provide overlapping or similar data at different quality levels. If necessary such federation will be artificially created. However,

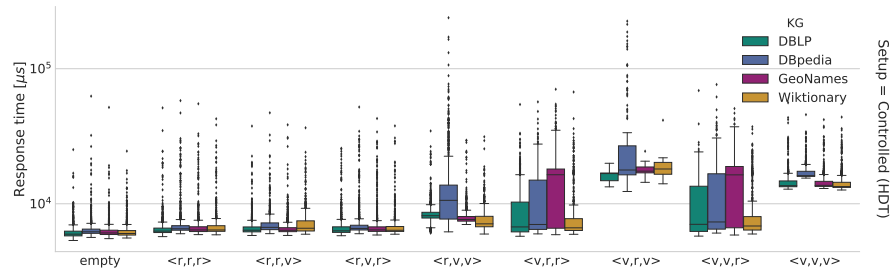


Fig. 2: Preliminary results: Boxplots for the response time of TPF server requesting different triple pattern types.

our goal is reusing, adapting and extending existing benchmarks (potentially also from relational databases) as far as possible.

Evaluation Metrics: We aim to consider several evaluation metrics including well-established metrics such as answer completeness and execution time to measure the performance of our solution. However, as we are especially interested in the quality of such query plans and the resulting answers, we will need to define a quality metric which allows assessing the quality of query plans considering several data quality metrics on both data source and query-level. Investigating these metrics will, therefore, provide an insight into the potential trade-off between answer completeness, execution time, and the quality of the query plans.

6 Preliminary Results

In previous research [11], we have investigated on hypothesis **H1** and especially the question of determining fine-grained statistics on the quality of RDF data sources. In our work on querying large knowledge graphs over Triple Pattern Fragments (TPFs) [11], we propose the *TPF Profiler* as a resource to retrieve performance statistics of TPF server. For a given TPF server URI and a sample size, the TPF Profiler collects performance measurements when submitting different requests to the server in three steps. First, the TPF Profiler randomly selects a sample of RDF triples from the RDF graph available at the server. Thereafter, a set of triple patterns is generated based on the sampled triples by replacing RDF terms with variables. Finally, the TPF Profiler requests all triple patterns generated in the previous step and measures the response time of the server. The resulting performance data may be analyzed with respect to various factors, such as the number and position of the variable in the triple pattern also referred to as triple pattern type. Figure 2 shows the measured response time from our empirical study for different Knowledge Graphs (KGs) according to the triple pattern type. The results indicate significant differences in response time for different triple pattern types. Hence, the results allow for deriving a triple pattern level quality assessment for the performance dimension which may be used to estimate the overall execution time of query plans.

Summarizing, our preliminary results show that the quality of data sources may vary on query-level, specifically on triple pattern level. Therefore, the results are a first step towards answering our first research question and may be used as a starting point for investigating the remaining hypotheses.

7 Conclusions and Lessons Learned

In this doctoral work, we aim to study the problem of quality-driven query plan optimization problem in the context of federated RDF data sources. We formulate three hypotheses to address the associated research questions and present an evaluation plan to investigate the validity of the hypotheses. Our preliminary results on **H1** provide insights into how query-specific data quality descriptions for data sources on the performance dimension may be assessed and potentially leveraged in query planning. However, it is still an open issue of how our empirical results can be integrated into a strategy for query planning.

Future work addressing **H1** will investigate further query-specific data source description for other quality dimensions as well as adapting existing data quality metrics to be efficiently measured via SPARQL endpoints and TPF servers. Thereafter, we want to address the remaining hypotheses **H2** and **H3**. Integral to this task is providing means to assess the quality of a query plan which takes various quality dimensions into account. Lastly, we need to determine heuristics suitable to explore the space of feasible query plans to find the best query plan.

Acknowledgements. I would like to thank my advisors Dr. Maribel Acosta and Prof. Dr. York Sure-Vetter for their support and valuable feedback.

References

1. Acosta, M., Hartig, O., Sequeda, J.: Federated RDF query processing. In: Sherif Sakr, A.Z. (ed.) *Encyclopedia of Big Data Technologies*. Springer (2018)
2. Acosta, M., Simperl, E., Flöck, F., Vidal, M.E.: Enhancing answer completeness of SPARQL queries via crowdsourcing. *Journal of Web Semantics* 45 (2017)
3. Acosta, M., Vidal, M.E., Lampo, T., Castillo, J., Ruckhaus, E.: ANAPSID: An Adaptive Query Processing Engine for SPARQL Endpoints. In: Aroyo, L., Welty, C., Alani, H., Taylor, J., Bernstein, A., Kagal, L., Noy, N., Blomqvist, E. (eds.) *The Semantic Web – ISWC 2011*, vol. 7031, pp. 18–34. Springer Berlin Heidelberg, Berlin, Heidelberg (2011)
4. Ben Ellef, M., Bellahsene, Z., Breslin, J.G., Demidova, E., Dietze, S., Szymański, J., Todorov, K.: RDF dataset profiling – a survey of features, methods, vocabularies and applications. *Semantic Web* 9(5), 677–705 (Aug 2018)
5. Darari, F., Nutt, W., Pirrò, G., Razniewski, S.: Completeness statements about RDF data sources and their use for query answering. In: *Advanced Information Systems Engineering*, vol. 7908, pp. 66–83. Springer (2013)
6. Endris, K.M., Galkin, M., Lytra, I., Mami, M.N., Vidal, M.E., Auer, S.: MULDER: Querying the Linked Data Web by bridging RDF molecule templates. In: Benslimane, D., Damiani, E., Grosky, W.I., Hameurlain, A., Sheth, A., Wagner, R.R. (eds.) *Database and Expert Systems Applications*, vol. 10438, pp. 3–18. Springer International Publishing, Cham (2017)

7. Färber, M., Bartscherer, F., Menne, C., Rettinger, A.: Linked data quality of DBpedia, Freebase, OpenCyc, Wikidata, and YAGO. *Semantic Web* 9(1) (2017)
8. Görlitz, O., Staab, S.: Splendid: Sparql endpoint federation exploiting VoID descriptions. In: *Proceedings of the Second International Conference on Consuming Linked Data-Volume 782*. pp. 13–24. CEUR-WS.org (2011)
9. Harth, A., Hose, K., Karnstedt, M., Polleres, A., Sattler, K.U., Umbrich, J.: Data summaries for on-demand queries over linked data. In: *Proceedings of the 19th international conference on World wide web - WWW '10*. p. 411. ACM Press, Raleigh, North Carolina, USA (2010)
10. Hartig, O.: Querying trust in RDF data with tSPARQL. In: Aroyo, L., Traverso, P., Ciravegna, F., Cimiano, P., Heath, T., Hyvönen, E., Mizoguchi, R., Oren, E., Sabou, M., Simperl, E. (eds.) *The Semantic Web: Research and Applications*, vol. 5554, pp. 5–20. Springer Berlin Heidelberg, Berlin, Heidelberg (2009)
11. Heling, L., Acosta, M., Maleshkova, M., Sure-Vetter, Y.: Querying large knowledge graphs over triple pattern fragments: An empirical study. In: *The Semantic Web - ISWC 2018 - 17th International Semantic Web Conference*, Monterey, CA, USA, October 8-12, 2018, *Proceedings, Part II*. pp. 86–102 (2018)
12. Hui, J., Li, L., Zhang, Z.: Integration of Big Data: A survey. In: Zhou, Q., Gan, Y., Jing, W., Song, X., Wang, Y., Lu, Z. (eds.) *Data Science*, vol. 901, pp. 101–121. Springer Singapore, Singapore (2018)
13. Ibaraki, T., Kameda, T.: On the optimal nesting order for computing N-relational joins. *ACM Transactions on Database Systems* 9(3), 482–502 (Aug 1984)
14. Lopes, N., Polleres, A., Straccia, U., Zimmermann, A.: AnQL: SPARQLing up annotated RDFS. In: Patel-Schneider, P.F., Pan, Y., Hitzler, P., Mika, P., Zhang, L., Pan, J.Z., Horrocks, I., Glimm, B. (eds.) *The Semantic Web – ISWC 2010*, vol. 6496, pp. 518–533. Springer Berlin Heidelberg (2010)
15. Naumann, F., Leser, U., Freytag, J.C.: Quality-driven integration of heterogenous information systems. In: *VLDB'99, Proceedings of 25th International Conference on Very Large Data Bases*, Edinburgh, Scotland, UK. pp. 447–458 (1999)
16. Neumann, T., Moerkotte, G.: Characteristic sets: Accurate cardinality estimation for RDF queries with multiple joins. In: *Data Engineering (ICDE), 2011 IEEE 27th International Conference on*. pp. 984–994 (Apr 2011)
17. Quilitz, B., Leser, U.: Querying distributed RDF data sources with SPARQL. In: Bechhofer, S., Hauswirth, M., Hoffmann, J., Koubarakis, M. (eds.) *The Semantic Web: Research and Applications*, vol. 5021, pp. 524–538. Springer Berlin Heidelberg, Berlin, Heidelberg (2008)
18. Saleem, M., Ngonga Ngomo, A.C.: HiBISCuS: Hypergraph-based source selection for SPARQL endpoint federation. In: *The Semantic Web: Trends and Challenges*, vol. 8465, pp. 176–191. Springer International Publishing, Cham (2014)
19. Schwarte, A., Haase, P., Hose, K., Schenkel, R., Schmidt, M.: Fedx: Optimization techniques for federated query processing on linked data. In: *International Semantic Web Conference*. pp. 601–616. Springer (2011)
20. Tsialiamanis, P., Sidirourgos, L., Fundulaki, I., Christophides, V., Boncz, P.: Heuristics-based query optimisation for SPARQL. In: *Proceedings of the 15th International Conference on Extending Database Technology - EDBT '12*. p. 324. ACM Press, Berlin, Germany (2012)
21. Wang, R.Y., Strong, D.M.: Beyond accuracy: What data quality means to data consumers. *Journal of Management Information Systems* 12(4), 5–33 (1996)
22. Zaveri, A., Rula, A., Maurino, A., Pietrobon, R., Lehmann, J., Auer, S.: Quality assessment for linked data: A survey. *Semantic Web* 7(1), 63–93 (2016)