

Enabling Fine-grained RDF Data Completeness Assessment

Fariz Darari*, Simon Razniewski, Radityo Eko Prasajo, and Werner Nutt

Free University of Bozen-Bolzano, Italy

*fariz.darari@stud-inf.unibz.it

Abstract. Nowadays, more and more RDF data is becoming available on the Semantic Web. While the Semantic Web is generally incomplete by nature, on certain topics, it already contains complete information and thus, queries may return all answers that exist in reality. In this paper we develop a technique to check query completeness based on RDF data annotated with completeness information, taking into account data-specific inferences that lead to an inference problem which is Π_2^F -complete. We then identify a practically relevant fragment of completeness information, suitable for crowdsourced, entity-centric RDF data sources such as Wikidata, for which we develop an indexing technique that allows to scale completeness reasoning to Wikidata-scale data sources. We verify the applicability of our framework using Wikidata and develop COOL-WD, a completeness tool for Wikidata, used to annotate Wikidata with completeness statements and reason about the completeness of query answers over Wikidata. The tool is available at <http://cool-wd.inf.unibz.it/>.

Keywords: RDF, data completeness, SPARQL, query completeness, Wikidata

1 Introduction

Over the Web, we are witnessing a growing amount of data available in RDF. The LOD Cloud¹ recorded that there were 1014 RDF data sources in 2014, covering various domains from life science to government. RDF follows the Open-World Assumption (OWA), assuming data is incomplete by default [1]. Yet, given such a large quantity of RDF data, one might wonder if it is complete for some topics. As an illustration, consider Wikidata, a crowdsourced KB with RDF support [2]. For data about the movie *Reservoir Dogs*, Wikidata is incomplete, as it is missing the fact that Michael Sottile was acting in the movie.² On the other hand, for data about *Apollo 11*, it is the case that Neil Armstrong, Buzz Aldrin, and Michael Collins, who are recorded as crew members on Wikidata, are indeed *all* the crew (see Figure 1).³ However, such completeness information is not recorded and thus it is left to the reader to decide whether some data on the Web is already complete.

¹ <http://lod-cloud.net/>

² By comparing the data at <https://www.wikidata.org/wiki/Q72962> with the complete information at <http://www.imdb.com/title/tt0105236/fullcredits>

³ <http://www.space.com/16758-apollo-11-first-moon-landing.html>



Fig. 1. Wikidata is actually complete for all the Apollo 11 crew

Nevertheless, the availability of explicit completeness information can benefit data access over RDF data sources, commonly done via SPARQL queries. To illustrate, suppose that in addition to the complete data of Apollo 11 crew, Wikidata is also complete for the children of the three astronauts. Consequently, a user asking for the children of Apollo 11 crew should obtain not only query answers, but also the information that the answers are complete.

Motivated by the above rationales, we argue that it is important to describe completeness of RDF data and provide a technique to check query completeness based on RDF data with its completeness information. Such a check is called *completeness entailment*. In previous work, Darari et al. [3] proposed a framework to describe completeness of RDF data and check query completeness based on completeness information. One fundamental limitation of this work is that the completeness check is agnostic of the content of the RDF data to which the completeness information is given, which results in weaker inferences. In the next section, we show that incorporating the content of RDF data may provide stronger inferences about query completeness. From the relational databases, Razniewski et al. [4], proposed wildcard-based completeness patterns to provide completeness information over databases. To check query completeness, they defined a pattern algebra, which works upon database tables enriched with completeness patterns. The work incorporated database instances in completeness check, which are conceptually similar to the content of RDF data. However, only a sound algorithm was provided for completeness check.

In this work, we make the following contributions:

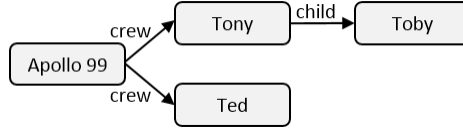
1. We provide a formalization of the completeness entailment problem for RDF data, and develop a sound and complete algorithm to solve the completeness entailment problem.
2. We identify a practically relevant fragment of completeness information suitable for crowdsourced, entity-centric RDF data sources like Wikidata, and develop an indexing technique to improve the feasibility of completeness entailment within the fragment.
3. We develop COOL-WD, a tool to manage completeness over Wikidata.

Our paper is structured as follows: Section 2 presents a motivating scenario. In Section 3, we provide a formalization to the completeness problem, followed

by Section 4 where we describe formal notions and a generic algorithm to check completeness entailment. Section 5 introduces a fragment of completeness information, suitable for crowdsourced, entity-centric RDF KBs like Wikidata, and presents an optimization technique for checking completeness entailment within this fragment. Section 6 reports our experimental evaluations. In Section 7, we describe COOL-WD. Related work is given in Section 8, whereas further discussion about our work is in Section 9. Section 10 concludes the paper and sketches future work.

2 Motivating Scenario

Let us consider a motivating scenario for the main problem of this work, that is, the check of query completeness based on RDF data with its completeness information. Consider an RDF graph about the crew of Apollo 99 (or for short, A99), a fictional space mission, and the children of the crew, as shown below.

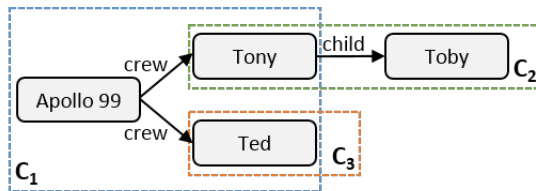


Consider now the query Q_0 asking for the crew of A99 and their children:

$$Q_0 = (W_0, P_0) = (\{ ?crew, ?child \}, \{ (a99, crew, ?crew), (?crew, child, ?child) \})$$

Evaluating Q_0 over the graph gives only one mapping result, where the crew is mapped to Tony and the child is mapped to Toby. Up until now, nothing can be said about the completeness of the query since: (i) there can be another crew member of A99 with a child; (ii) Tony may have another child; or (iii) Ted may have a child.

Let us now consider the same graph as before, now enriched with completeness information, as shown below.



Informally, the above figure contains three completeness statements: C_1 , which states that the graph contains all crew members of A99; C_2 , which states the graph contains all Tony's children; and C_3 , which states the graph contains all Ted's children (i.e., Ted has no children). With the addition of completeness information, let us see whether we may answer our query completely.

First, since we know that all the crew of A99 are Tony and Ted, the query Q_0 then becomes equivalent to the following two queries:

- $Q_1 = (W_1, P_1) = (\{ ?child \}, \{ (a99, crew, tony), (tony, child, ?child) \})$
- $Q_2 = (W_2, P_2) = (\{ ?child \}, \{ (a99, crew, ted), (ted, child, ?child) \})$

where the variable $?crew$ is instantiated with Tony and Ted, respectively.

Moreover, for Q_2 according to our graph and completeness information, Ted has no children. Thus, there is no way that Q_2 will return an answer, so Q_2 can be safely removed. Now, only Q_1 is left. Again, from our graph and completeness information, we know that Toby is the only child of Tony. Thus, Q_1 in turn is equivalent to the following boolean query:

$$Q_3 = (W_3, P_3) = (\{ \}, \{ (a99, crew, tony), (tony, child, toby) \})$$

with the variable $?crew$ is instantiated to Tony and $?child$ to Toby. However, our graph is complete for Q_3 as it contains the whole body of Q_3 . Since from our reasoning the query Q_3 is equivalent to our original query Q_0 , we conclude that our graph with its completeness information can guarantee the completeness of Q_0 , that is, Toby is the only child of Tony, the only crew member of A99 having a child.

Note that using the data-agnostic approach from [3], it is not possible to derive the same conclusion. Without looking at the actual graph, we cannot conclude that Ted and Tony are all the crew members of Apollo 99. Consequently, just having the children of Tony and Ted complete does not help in reasoning about Apollo 99. In the rest of the paper, we discuss how the intuitive, data-specific reasoning from above can be formalized.

3 Formal Framework

In this section, we remind the reader of RDF and SPARQL, and provide formalization to our completeness problem.

RDF and SPARQL. Assume three pairwise disjoint infinite sets I (*IRIs*), L (*literals*), and V (*variables*). A tuple $(s, p, o) \in I \times I \times (I \cup L)$ is called a triple, while a finite set of triples is called an RDF graph.

The standard query language for RDF graphs is SPARQL [5]. At the core of SPARQL lies triple patterns, which are like triples, but also variables are allowed in each position. In this work, we focus on the conjunctive fragment of SPARQL, which uses sets of triple patterns, called basic graph patterns (BGPs). Evaluating a BGP P over G gives the set of mappings $\llbracket P \rrbracket_G = \{ \mu \mid \mu P \subseteq G \text{ and } \text{dom}(\mu) = \text{var}(P) \}$. Over P , we define a *freeze mapping* \tilde{id} that maps each variable $?v$ in P to a fresh IRI \tilde{v} . From such a mapping, we construct *the prototypical graph* $\tilde{P} := \tilde{id}P$ that encodes any possible graph that can satisfy the BGP. A query $Q = (W, P)$ projects the evaluation results of a BGP P to a set W of variables. Moreover, a **CONSTRUCT** query has the abstract form $(\text{CONSTRUCT } P_1 P_2)$ where both P_1 and P_2 are BGPs. Evaluating a **CONSTRUCT** query over G results in a graph where P_1 is instantiated with all the mappings in $\llbracket P_2 \rrbracket_G$.

Completeness Statements. A completeness statement describes which parts of an RDF graph are complete. We adopt the definition of completeness statements in [3].

Definition 1 (Completeness Statement). *A completeness statement C is defined as $\text{Compl}(P_C)$ where P_C is a non-empty BGP.*

Example 1. The completeness statements in our motivating scenario are as follows: $C_1 = \text{Compl}(\{(a99, \text{crew}, ?c)\})$, $C_2 = \text{Compl}(\{(tony, \text{child}, ?c)\})$, and $C_3 = \text{Compl}(\{(ted, \text{child}, ?c)\})$.

To serialize completeness statements in RDF, we refer the reader to [3]. Now, let us define the semantics of completeness statements. First, we associate the CONSTRUCT query $Q_C = (\text{CONSTRUCT } P_C \text{ } P_C)$ to each statement C . From now on, we fix a graph G upon which we describe its completeness. Given a graph $G' \supseteq G$, we call (G, G') an *extension pair*. In general, with no completeness statement, every extension pair is a valid extension pair, that is, G' is a possible state of the ideal world where all the information is complete. For instance, without completeness statement, in the motivating scenario, all of the following would be valid extensions: That there are more crew members of A99; that Tony has more children; and that Ted has children. Completeness statements restrict the valid extensions of a graph.

Definition 2 (Valid Extension Pairs). *Let G be a graph and C a completeness statement. We say that an extension pair (G, G') is valid wrt. C , written $(G, G') \models C$, if $\llbracket Q_C \rrbracket_{G'} \subseteq G$.*

The above definition naturally extends to sets of completeness statements. Over a set \mathcal{C} of completeness statements and a graph G , we define *the transfer operator* $T_{\mathcal{C}}(G) = \bigcup_{C \in \mathcal{C}} \llbracket Q_C \rrbracket_G$. We have the following characterization: for all extension pairs (G, G') , it is the case that $(G, G') \models \mathcal{C}$ iff $T_{\mathcal{C}}(G') \subseteq G$.

Query Completeness. We write $\text{Compl}(Q)$ to denote that a query Q is complete. Over an extension pair, a query is complete iff it returns the same results over both the graphs of the extension pair.

Definition 3 (Query Completeness). *Let (G, G') be an extension pair and Q be a query. We define that: $(G, G') \models \text{Compl}(Q)$ iff $\llbracket Q \rrbracket_{G'} = \llbracket Q \rrbracket_G$.⁴*

Completeness Entailment We now define the main problem of our work, the completeness entailment.

Definition 4 (Completeness Entailment). *Given a set \mathcal{C} of completeness statements, a graph G , and a query Q , we define that \mathcal{C} and G entail the completeness of Q , written as $\mathcal{C}, G \models \text{Compl}(Q)$, if for all extension pairs $(G, G') \models \mathcal{C}$, it holds that $(G, G') \models \text{Compl}(Q)$.*

⁴ Since in this work we focus on conjunctive queries which are monotonic, the direction $\llbracket Q \rrbracket_{G'} \supseteq \llbracket Q \rrbracket_G$ comes for free.

In our motivating scenario, we have seen that the graph about the crew of A99 and the completeness statements there entail the completeness of the query Q_0 asking for the crew of A99 and their children.

In this work, we assume bag semantics for query evaluation, which is the default of SPARQL.⁵ Consequently, this allows us to focus on the BGPs used in the body of conjunctive queries for completeness entailment.

4 Checking Completeness Entailment

In this section, we present an algorithm for performing the completeness check as demonstrated in our motivating scenario.

4.1 Preliminaries

Before presenting the algorithm, we introduce important notions.

First, we need to have a notion for a BGP with a stored mapping from variable instantiations. Let P be a BGP and μ be a mapping such that $\text{dom}(\mu) \cap \text{var}(P) = \emptyset$. We define the pair (P, μ) as a *partially mapped BGP*, which is a BGP with a stored mapping. Over a graph G , the evaluation of (P, μ) is defined as $\llbracket (P, \mu) \rrbracket_G = \{ \mu \cup \nu \mid \nu \in \llbracket P \rrbracket_G \}$. It is easy to see that $P \equiv (P, \emptyset)$. Furthermore, we define the evaluation of a set of partially mapped BGPs over a graph G as the union of evaluating each of them over G .

Example 2. Consider our motivating scenario. Over the BGP P_0 of the query Q_0 , instantiating the variable $?crew$ to *tony* results in the BGP P_1 of the query Q_1 . Pairing P_1 with this instantiation gives the partially mapped BGP $(P_1, \{?crew \mapsto \text{tony}\})$. Moreover, it is the case that $\llbracket (P_1, \{?crew \mapsto \text{tony}\}) \rrbracket_G = \{ \{?crew \mapsto \text{tony}, ?child \mapsto \text{toby}\} \}$.

Next, we would like to formalize the equivalence between partially mapped BGPs wrt. a set \mathcal{C} of completeness statements and a graph G .

Definition 5 (Equivalence under \mathcal{C} and G). *Let (P, μ) and (P', ν) be partially mapped BGPs, \mathcal{C} be a set of completeness statements, and G be a graph. We define that (P, μ) is equivalent to (P', ν) wrt. \mathcal{C} and G , written $(P, \mu) \equiv_{\mathcal{C}, G} (P', \nu)$, if for all $(G, G') \models \mathcal{C}$, it holds that $\llbracket (P, \mu) \rrbracket_{G'} = \llbracket (P', \nu) \rrbracket_{G'}$.*

The above definition naturally extends to sets of partially mapped BGPs.

Example 3. Consider all the queries in our motivating scenario. It is the case that $\{ (P_0, \emptyset) \} \equiv_{\mathcal{C}, G} \{ (P_1, \{?crew \mapsto \text{tony}\}), (P_2, \{?crew \mapsto \text{ted}\}) \} \equiv_{\mathcal{C}, G} \{ (P_3, \{?crew \mapsto \text{tony}, ?child \mapsto \text{toby}\}) \}$.

Next, we would like to figure out which parts of a query contain variables that can be instantiated completely. For this reason, we define

$$\text{cruc}_{\mathcal{C}, G}(P) = P \cap \tilde{id}^{-1}(T_{\mathcal{C}}(\tilde{P} \cup G))$$

⁵ <http://www.w3.org/TR/sparql11-query/>

as the crucial part of P wrt. \mathcal{C} and G . It is the case that $\mathcal{C}, G \models \text{Compl}(\text{cruc}_{\mathcal{C},G}(P))$, that is, we are complete for the crucial part. Later on, we will see that the crucial part is used to ‘guide’ the instantiation process during the completeness entailment check.

Example 4. Consider the query $Q_0 = (W_0, P_0)$ in our motivating scenario. We have that $\text{cruc}_{\mathcal{C},G}(P_0) = P_0 \cap \tilde{id}^{-1}(T_{\mathcal{C}}(\tilde{P}_0 \cup G)) = \{(a99, \text{crew}, ?\text{crew})\}$ with $\tilde{id} = \{?\text{crew} \mapsto \widetilde{\text{crew}}, ?\text{child} \mapsto \widetilde{\text{child}}\}$. Consequently, we can have a complete instantiation of the crew of A99.

The operator below implements the instantiations of a partially mapped BGP wrt. its crucial part.

Definition 6 (Equivalent Partial Grounding). *Let \mathcal{C} be a set of completeness statements, G be a graph, and (P, ν) be a partially mapped BGP. We define the operator equivalent partial grounding:*

$$\text{epg}((P, \nu), \mathcal{C}, G) = \{(\mu P, \nu \cup \mu) \mid \mu \in \llbracket \text{cruc}_{\mathcal{C},G}(P) \rrbracket_G\}.$$

The following lemma shows that such instantiations produce a set of partially mapped BGPs equivalent to the original partially mapped BGP, hence the name equivalent partial grounding. The lemma holds since the instantiation is done over the crucial part, which is complete wrt. \mathcal{C} and G .

Lemma 1 (Equivalent Partial Grounding). *Let \mathcal{C} be a set of completeness statements, G be a graph, and (P, ν) be a partially mapped BGP. We have that*

$$\{(P, \nu)\} \equiv_{\mathcal{C},G} \text{epg}((P, \nu), \mathcal{C}, G).$$

Example 5. Consider our motivating scenario. We have that:

- $\text{epg}((P_2, \{?\text{crew} \mapsto \text{ted}\}), \mathcal{C}, G) = \emptyset$
- $\text{epg}((P_3, \{?\text{crew} \mapsto \text{tony}, ?\text{child} \mapsto \text{toby}\}), \mathcal{C}, G) = \{(P_3, \{?\text{crew} \mapsto \text{tony}, ?\text{child} \mapsto \text{toby}\})\}$
- $\text{epg}((P_0, \emptyset), \mathcal{C}, G) = \{(P_1, \{?\text{crew} \mapsto \text{tony}\}), (P_2, \{?\text{crew} \mapsto \text{ted}\})\}$

Generalizing from the example above, there are three cases of $\text{epg}((P, \nu), \mathcal{C}, G)$:

- If $\llbracket \text{cruc}_{\mathcal{C},G}(P) \rrbracket_G = \emptyset$, it returns an empty set.
- If $\llbracket \text{cruc}_{\mathcal{C},G}(P) \rrbracket_G = \{\emptyset\}$, it returns $\{(P, \nu)\}$.
- Otherwise, it returns a non-empty set of partially mapped BGPs where some variables in P are instantiated.

From these three cases and the finite number of triple patterns with variables of a BGP, it holds that the repeated applications of the epg operator, with the first and second cases above as the base cases, are terminating. Note that the difference between these two base cases is on the effect of their corresponding epg operations, as illustrated in Example 5: for the first case, the epg operation returns an empty set, whereas for the second case, it returns back the input partially mapped BGP.

We define that a partially mapped BGP (P, ν) is *saturated* wrt. \mathcal{C} and G , if $epg((P, \nu), \mathcal{C}, G) = \{(P, \nu)\}$, that is, if the second case above applies. Note that the notion of saturation is independent from the mapping in a partially mapped BGP: given a mapping ν , a partially mapped BGP (P, ν) is saturated wrt. \mathcal{C} and G iff (P, ν') is saturated wrt. \mathcal{C} and G for any mapping ν' . Thus, wrt. \mathcal{C} and G we say that a BGP P is saturated if (P, \emptyset) is saturated.

The completeness checking of saturated BGPs is straightforward as we only need to check if they are contained in the graph G .

Proposition 1 (Completeness Entailment of Saturated BGPs). *Let P be a BGP, \mathcal{C} be a set of completeness statements, and G be a graph. Suppose P is saturated wrt. \mathcal{C} and G . Then, it is the case that: $\mathcal{C}, G \models \text{Compl}(P)$ iff $\tilde{P} \subseteq G$.*

Based on the above notions, we are ready to provide an algorithm to check completeness entailment. The next subsection gives the algorithm.

4.2 Algorithm for Checking Completeness Entailment

Now we introduce an algorithm to compute all saturated, equivalent partial grounding results of a BGP wrt. \mathcal{C} and G . Following from Proposition 1, we can then check whether all the resulting saturated BGPs are contained in the graph G to see if the completeness entailment holds.

ALGORITHM 1: $\text{sat}(P_{\text{orig}}, \mathcal{C}, G)$

Input: A BGP P_{orig} , a set \mathcal{C} of completeness statements, a graph G

Output: A set Ω of mappings

```

1  $\mathbf{P}_{\text{working}} \leftarrow \{(P_{\text{orig}}, \emptyset)\}$ 
2  $\Omega \leftarrow \emptyset$ 
3 while  $\mathbf{P}_{\text{working}} \neq \emptyset$  do
4    $(P, \nu) \leftarrow \mathbf{takeOne}(\mathbf{P}_{\text{working}})$ 
5    $\mathbf{P}_{\text{equiv}} \leftarrow \text{epg}((P, \nu), \mathcal{C}, G)$ 
6   if  $\mathbf{P}_{\text{equiv}} = \{(P, \nu)\}$  then
7      $\Omega \leftarrow \Omega \cup \nu$ 
8   else
9      $\mathbf{P}_{\text{working}} \leftarrow \mathbf{P}_{\text{working}} \cup \mathbf{P}_{\text{equiv}}$ 
10  end
11 end
12 return  $\Omega$ 

```

Consider a BGP P_{orig} , a set \mathcal{C} of completeness statements, and a graph G . The algorithm works as follows: First, we transform our original BGP P_{orig} into its equivalent partially mapped BGP $(P_{\text{orig}}, \emptyset)$ and put it in $\mathbf{P}_{\text{working}}$. Then, in each iteration of the while loop, we take and remove a partially mapped BGP (P, ν) from $\mathbf{P}_{\text{working}}$ via the method **takeOne**. Afterwards, we compute $epg((P, \nu), \mathcal{C}, G)$. As discussed above there might be three result cases here: (i) If $epg((P, \nu), \mathcal{C}, G) = \emptyset$, then simply we remove (P, ν) and will not consider it anymore in the later iteration; (ii) If $epg((P, \nu), \mathcal{C}, G) = \{(P, \nu)\}$, that is, (P, ν) is saturated, then we collect the mapping ν to the set Ω ; and (iii) otherwise,

we add to $\mathbf{P}_{working}$ a set of partially mapped BGPs instantiated from (P, ν) . We keep iterating until $\mathbf{P}_{working} = \emptyset$, and finally return the set Ω .

The following proposition follows from the construction of the above algorithm and Lemma 1.

Proposition 2. *Given a BGP P , a set \mathcal{C} of completeness statements, and a graph G , the following properties hold:*

- For all $\mu \in \text{sat}(P, \mathcal{C}, G)$, it is the case that μP is saturated wrt. \mathcal{C} and G .
- It holds that $\{(P, \emptyset)\} \equiv_{\mathcal{C}, G} \{(\mu P, \mu) \mid \mu \in \text{sat}(P, \mathcal{C}, G)\}$.

From the above proposition, we can derive the following theorem, which shows the soundness and completeness of the algorithm to check completeness entailment.

Theorem 1 (Completeness Entailment Check). *Let P be a BGP, \mathcal{C} be a set of completeness statements, and G be a graph. It holds that*

$$\mathcal{C}, G \models \text{Compl}(P) \quad \text{iff} \quad \text{for all } \mu \in \text{sat}(P, \mathcal{C}, G) . \widetilde{\mu P} \subseteq G.$$

Example 6. Consider our motivating scenario. We have that $\text{sat}(P_0, \mathcal{C}, G) = \{\{?crew \mapsto \text{tony}, ?child \mapsto \text{toby}\}\}$. It is the case that for all $\mu \in \text{sat}(P_0, \mathcal{C}, G)$, it holds that $\widetilde{\mu P_0} \subseteq G$. Thus, by Theorem 1 the entailment $\mathcal{C}, G \models \text{Compl}(P)$ holds.

By reduction from validity of $\forall\exists$ SAT formula, one can show that the complexity of the completeness entailment is Π_2^P -complete.

Corollary 1 (Complexity of Completeness Check). *Deciding whether the entailment $\mathcal{C}, G \models \text{Compl}(P)$ holds, given a set \mathcal{C} of completeness statements, a graph G , and a BGP P , is Π_2^P -complete.*

In what follows, we provide optimization techniques for the algorithm, which work for generic cases of completeness entailment.

Early failure detection. In our algorithm, the containment checks for saturated BGPs are done at the end. Indeed, if there is a single saturated BGP not contained in the graph, we cannot guarantee query completeness. Thus, instead of having to collect all saturated BGPs and then check the containment later on, we can improve the performance of the algorithm by performing the containment check right after the saturation check (Line 6 of the algorithm). So, as soon as there is a failure in the containment check, we stop the loop and conclude that the completeness entailment does not hold.

Completeness skip. Recall the definition of the operator $epg((P, \nu), \mathcal{C}, G) = \{(\mu P, \nu \cup \mu) \mid \mu \in \llbracket crucc_{\mathcal{C}, G}(P) \rrbracket_G\}$, which relies on the `cruc` operator. Now, suppose that $crucc_{\mathcal{C}, G}(P) = P$, that is, we are complete for the whole part of the BGP P . Thus, we actually do not have to instantiate P in the `epg` operator, since we know that the instantiation results are contained in G as the consequence of it being complete wrt. \mathcal{C} and G . In conclusion, whenever $crucc_{\mathcal{C}, G}(P) = P$, we just remove (P, ν) from $\mathbf{P}_{working}$ and thus skip its instantiations.

Despite these optimizations, for a large number of completeness statements, the completeness entailment check may take long. In the next section, we identify a practically relevant fragment of completeness statements, for which we develop an indexing technique to make the entailment check feasible.

5 A Practical Fragment of Completeness Statements

This section identifies *SP-statements*, a fragment of completeness statements possessing several properties that are suitable to be used in practice. In the next sections, we show by experimental evaluations the feasibility of this fragment with an indexing technique we describe below, and demonstrate a completeness tool for Wikidata using the fragment.

5.1 SP-Statements

An *SP-statement* $Compl(\{(s, p, ?v)\})$ is a completeness statement with only one triple pattern in the BGP of the statement, where the subject and the predicate are IRIs, and the object is a variable.⁶ In our motivating scenario, all the completeness statements are in fact SP-statements. The statements possess the following properties, which are suitable for practical use:

- Having a simple structure, completeness statements within this fragment are easy to create and to be read. Thus, they are suitable for *crowdsourced* KBs, where humans are involved.
- An SP-statement denotes the completeness of all the property values of the entity which is the subject of the statement. This fits *entity-centric* KBs like Wikidata, where data is organized into entities (i.e., each entity has its own data page).
- Despite their simplicity, SP-statements can be used to guarantee the completeness of more complex queries such as queries whose length is greater than one (as illustrated by our motivating scenario).

5.2 SP-Indexing

We describe here how to optimize completeness entailment check with SP-statements. Recall our generic algorithm to check completeness entailment: In

⁶ We do not allow the subject to be a variable as it is not practically reasonable (e.g., complete for all the entities and values of predicate `child`).

the *cruc* operator within the *epg* operator (Line 5 of Algorithm 1), we have to compute $T_{\mathcal{C}}(\tilde{P} \cup G)$, that is, evaluate all **CONSTRUCT** queries of the completeness statements in \mathcal{C} over the graph $\tilde{P} \cup G$. This may be problematic if there are a large number of completeness statements in \mathcal{C} . Thus, we want to avoid such costly $T_{\mathcal{C}}$ applications. Given that completeness statements are of the form SP-statements, we may instead look for the statements having the same subject and predicate of the triple patterns in the BGP. The crucial part of the BGP P wrt. \mathcal{C} and G are the triple patterns with the matching subject and predicate of the completeness statements.

Proposition 3. *Given a BGP P , a graph G , and a set \mathcal{C} of SP-statements, it is the case that*
 $cruc_{\mathcal{C},G}(P) = \{ (s, p, o) \in P \mid \text{there exists a statement } Compl(\{ (s, p, ?v) \}) \in \mathcal{C} \}$.

From the above proposition, to get the crucial part, we only have to find an SP-statement with the same subject and predicate for each triple pattern of the BGP. In practice, we can facilitate this search using a standard hashmap, providing constant-time performance, also for other basic operations such as **add** and **delete**. The hashmap provides a mapping from the concatenation of the subject and the predicate of a statement to the statement itself. To illustrate, the hashmap of the completeness statements in our motivating scenario is as follows: $\{ a99-crew \mapsto C_1, tony-child \mapsto C_2, ted-child \mapsto C_3 \}$.

6 Experimental Evaluation

Now that we have an indexing technique for SP-statements, we want to see the performance of completeness check. To do so, we perform experimental evaluations with a realistic scenario, where we compare the runtime of completeness entailment when query completeness can be guaranteed (i.e., the success case), completeness entailment when query completeness cannot be guaranteed (i.e., the failure case), and query evaluation.

Experimental Setup. Our reasoning algorithm and indexing modules are implemented in Java using the Apache Jena library.⁷ We use Jena-TDB as the triple store of our experiment. The SP-indexing is implemented using the standard Java **hashmap**, where the keys are strings, constructed from the concatenation of the subject and predicate of completeness statements, and the values are Java objects representing completeness statements. All experiments are done on a standard laptop with a 2.4 GHz Intel Core i5 and 8 GB of memory.

To perform the experiment, we need three ingredients: a graph, completeness statements, and queries. For the graph, we use the direct-statement fragment of the Wikidata graph, which does not include qualifiers nor references and consists of 100 mio triples.⁸ The completeness statements and queries of this experiment are constructed based on the following pattern queries:

⁷ <https://jena.apache.org/>

⁸ http://tools.wmflabs.org/wikidata-exports/rdf/index.php?content=dump_download.php&dump=20151130

1. Give all mothers of mothers of mothers.
 $P_1 = \{ (?v, P25, ?w), (?w, P25, ?x), (?x, P25, ?y) \}$
2. Give the crew of a thing, the astronaut missions of that crew, and the operator of the missions.
 $P_2 = \{ (?v, P1029, ?w), (?w, P450, ?x), (?x, P137, ?y) \}$
3. Give the administrative divisions of a thing, the administrative divisions of those divisions, and their area.
 $P_3 = \{ (?v, P150, ?w), (?w, P150, ?x), (?x, P2046, ?y) \}$

To generate queries, we simply evaluate each pattern query over the graph, and instantiate the variable $?v$ of each pattern query with the corresponding mappings from the evaluation. We record 5200 queries instantiated from P_1 , 57 queries from P_2 , and 475 queries from P_3 . Each pattern query has a different average number of query results: the instantiations of P_1 give 1 result, those of P_2 give 4 results, and those of P_3 give 108 results on average.

To generate completeness statements, from each generated query, we iteratively evaluate each triple pattern from left to right, and construct SP-statements from the instantiated subject and the predicate of the triple patterns. This way, we guarantee that all the queries can be answered completely. We generate in total around 1.7 mio statements, with 30072 statements for P_1 , 484 statements for P_2 , and 1682263 statements for P_3 . Such a large number of completeness statements would make completeness checks without indexing very slow: Performing just a single application of the T_C operator with all these statements, which occurs in the `cruc` operator of the algorithm without SP-indexing, took about 20 minutes. Note that in a completeness check, there might be many T_C applications.

Now we describe how to observe the behavior when queries cannot be guaranteed to be complete, that is, the failure case. In this case, we drop randomly 20% of the completeness statements for each pattern query. To make up the statements we drop, we add dummy statements with the number equal to the number of dropped statements. This way, we ensure the same number of completeness statements for both the success and failure case.

For each query pattern, we measure the runtime of completeness check for both the success case and the failure case, and then query evaluation for the success case.⁹ We take 40 sample queries for each pattern query, repeat each run 10 times, and report the median of these runs.

Experimental Results. The experimental results are shown in Figure 2. Note that the runtime is in log scale. We can see that in all cases, the runtime increases with the first pattern query having the lowest runtime, and the third pattern query having the highest runtime. This is likely due to the increased number of query results. We observe that in all pattern queries, completeness check when queries are guaranteed to be complete is slower than those whose completeness cannot be guaranteed. We suspect that this is because in the former case, variable

⁹ We do not measure query evaluation time for failure case since query evaluation is independent of the completeness of the query.

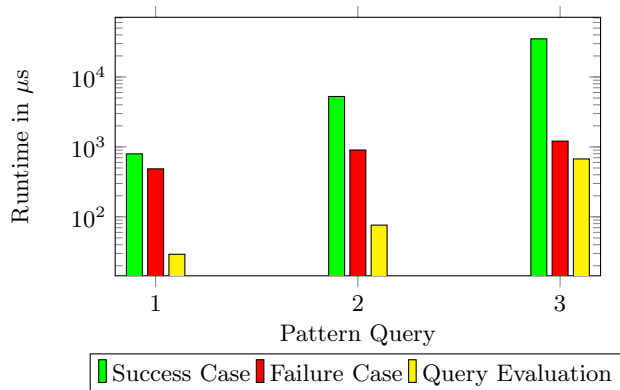


Fig. 2. Experiment Results of Completeness Entailment

instantiations have to be performed much more than in the latter case. In the latter case, as also described in Subsection 4.2, as soon as we find a saturated BGP not contained in the graph, we stop the loop in the algorithm and return `false`, meaning that the query completeness cannot be guaranteed.

In absolute scale, completeness check runs relatively fast, with $796 \mu s$ for P_1 , $5264 \mu s$ for P_2 , and $35130 \mu s$ for P_3 in success case; and $485 \mu s$ for P_1 , $903 \mu s$ for P_2 , and $1209 \mu s$ for P_3 in failure case. Note that as mentioned before, completeness check without indexing is not feasible at all here, as there are a large number of completeness statements, making the T_C application very slow (i.e., 20 minutes for a single application). For all pattern queries, however, query evaluation runs faster than completeness checking. This is because completeness checking may involve several query evaluations during the instantiation process with the `epg` operator.

To conclude, we have observed that completeness checking with a large number of SP-statements can be done reasonably fast, even for large datasets, by the employment of indexing. Also, we observe a clear positive correlation between the number of query results and the runtime of completeness checking. Last, performing completeness check when a query is complete is slower than that when a query cannot be guaranteed to be complete.

7 COOL-WD: A Completeness Tool for Wikidata

In this section, we introduce COOL-WD, a **CO**mpleteness to**OL** for **WikiData**. The tool implements our completeness framework with SP-statements and focuses to provide completeness information for direct statements of Wikidata. While our implementation is based on Apache Jena, our approach can be applied also via other Semantic Web frameworks like Sesame.¹⁰ Our tool is inspired by real, natural language completeness statements on Wikipedia, where com-

¹⁰ <http://rdf4j.org/>

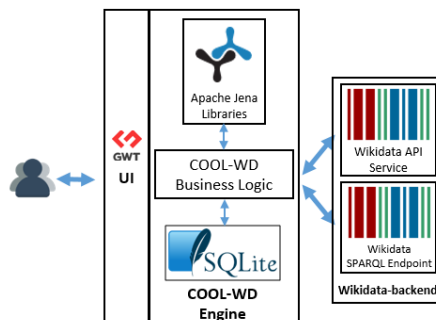


Fig. 3. COOL-WD Architecture

pleteness statements are given in a crowdsourced way.¹¹ The tool is available at <http://cool-wd.inf.unibz.it/>.

7.1 System Architecture

As shown in Figure 3, COOL-WD consists of three main components: user interface (UI), COOL-WD engine, and Wikidata-backend.

The first component is the UI, developed using GWT.¹² The UI provides the front-end interface for COOL-WD users, enabling them to search for Wikidata entities, look at facts about them enriched with completeness information, add/remove completeness statements, and check the completeness of a Wikidata query.

The second component is the engine, responsible for storing completeness statements using SQLite and performing completeness checks. We use optimization techniques as described in Subsection 4.2 and SP-indexing as described in Section 5 to improve the performance of completeness checks.

The last component is the Wikidata-backend. It consists of two subcomponents: Wikidata API and Wikidata SPARQL endpoint. The API is used for the suggestions feature in searching for Wikidata entities, while the Wikidata SPARQL endpoint serves as the source of Wikidata facts to which completeness statements are given, and of query evaluation.

7.2 Tool Usage

Here, we describe how one can use COOL-WD. From the landing page, the user is provided with a search bar for Wikidata entities. The search bar features auto-complete search suggestions, matching user keywords with the English labels of Wikidata entities. Clicking on a search suggestion gives the users the entity page, consisting of Wikidata facts about the entity with its completeness information.

¹¹ https://en.wikipedia.org/wiki/Template:Complete_list

¹² <http://www.gwtproject.org/>

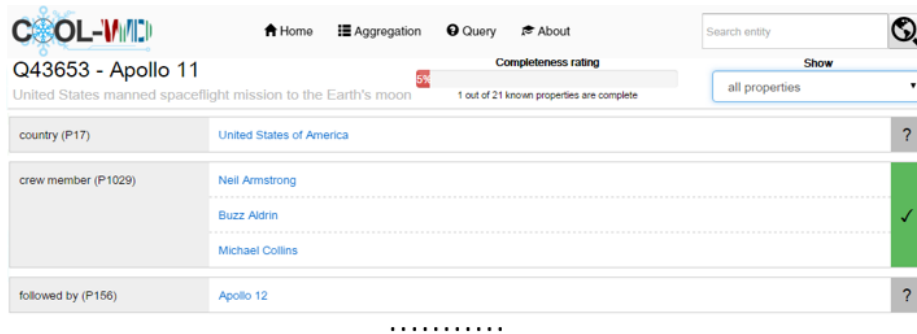


Fig. 4. The COOL-WD page of Apollo 11. Complete property values are with checkmarks.

An example is shown in Figure 4, which is the Apollo 11 page with the complete crew. Complete properties are distinguished by the checkmark symbol. To add a completeness statement, a user simply clicks a question mark next to the respective properties of an entity. Additionally, it is possible to add provenance information about authors, timestamps, and references of the statement. Suppose the user would also like to add completeness statements for the astronaut missions of Neil Armstrong, Buzz Aldrin, and Michael Collins. Therefore, she may perform an analogous operation: go to the entity pages, and click the question mark next to the respective properties. We also have a feature to see all stored completeness statements over Wikidata filtered by properties on the aggregation page.

If a user would like to evaluate a query and check its completeness, she has to go to the query page. Suppose she wants to know the crew of Apollo 11 and their astronaut missions. The user then specifies her query, and executes it. Instead of having only query answers, she can also see the completeness information of the answers.

8 Related Work

Data completeness concerns the breadth, depth, and scope of information [6]. In the relational databases, Motro [7] and Levy [8] were among the first to investigate data completeness. Motro developed a sound technique to check query completeness based on database views, while Levy introduced the notion of local completeness statements to denote which parts of a database are complete. Razniewski and Nutt [9] further extended their results by reducing completeness reasoning to containment checking, for which many algorithms are known, and characterizing the complexity of reasoning for different classes of queries. In terms of their terminology, our completeness entailment problem is one of QC-QC entailment under bag semantics, for which so far it was only known that it is in Π_3^P [10]. In [4], Razniewski et al. proposed completeness patterns and defined

a pattern algebra to check the completeness of queries. The work incorporated database instances, yet provided only a sound algorithm for completeness check.

We now move on to the Semantic Web. Fürber and Hepp [11] distinguished three types of completeness: ontology completeness, concerning which ontology classes and properties are represented; population completeness, referring to whether all objects of the real-world are represented; and property completeness, measuring the missing values of a specific property. In our work, SP-statements can be used to state the property completeness of an entity. Mendes et al. [12] proposed Sieve, a framework for expressing quality assessment and fusion methods, where completeness is also considered. With Sieve, users can specify how to compute quality scores and express a quality preference specifying which characteristics of data indicate higher quality. In the context of crowdsourcing, Chu et al. [13] developed KATARA, a hybrid data cleaning system, which not only cleans data, but may also add new facts to increase the completeness of the KB; whereas Acosta et al. [14] developed HARE, a hybrid SPARQL engine to enhance answer completeness.

Galárraga et al. [15] proposed a rule mining system that is able to operate under the Open-World Assumption (OWA) by simulating negative examples using the Partial Completeness Assumption (PCA). The PCA assumes that if the dataset knows some r -attribute of x , then it knows all r -attributes of x . This heuristic was also employed by Dong et al. [16] to develop Knowledge Vault, a Web-scale system for probabilistic knowledge fusion. In their paper, they used the term Local Closed-World Assumption (LCWA).

9 Discussion

We discuss here various aspects of our work: sources of completeness statements, completeness statements with provenance, and no-value information.

Sources of Completeness Statements. As demonstrated by COOL-WD, one way to provide completeness statements is via crowdsourcing. For domain-specific data like biology and archeology, domain experts may be a suitable source of completeness statements. An automated way to add completeness statements can also be leveraged by using NLP techniques to extract natural language completeness statements already available on the Web: around 13000 Wikipedia pages contain the keywords “complete list of” and “list is complete”, while IMDb provides complete cast information with the keywords “verified as complete” for some movies like Reservoir Dogs.¹³

Completeness Statements with Provenance. Just as data can be wrong, completeness statements can be wrong, too. Moreover, as data may change over time, completeness statements can be out-of-date. As a possible solution, one can add provenance information. Adding information about the author and reference of completeness statements may be useful to check the correctness of the

¹³ <http://www.imdb.com/title/tt0105236/fullcredits>

statements, while attaching timestamps would provide timeliness information to the statements.

No-Value Information. Completeness statements can also be used to represent the non-existence of information. For example, in our motivating scenario, there is the completeness statement about the children of Ted with no corresponding data in the graph. In this case, we basically say that Ted has no children. As a consequence of having no-value information, we can be complete for queries despite having the empty answer. Such a feature is similar to that proposed in [17]. The only difference is that here we need to pair completeness statements with a graph that has no corresponding data captured by the statements, while in that work, no-value statements are used to directly say that some parts of data do not exist.

10 Conclusions and Future Work

The availability of an enormous amount of RDF data calls for better data quality management. In this work, we focus on the data quality aspect of completeness. We develop a technique to check query completeness based on RDF data with its completeness information. To increase the practical benefits of our framework, we identify a practically relevant fragment of completeness information upon which an indexing can be implemented to optimize completeness check, and develop COOL-WD, a completeness management tool for Wikidata.

For future work, we would like to investigate indexing techniques for more general cases. One challenge here is that how to index the arbitrary structure of completeness statements. Another plan is to develop a technique to extract completeness statements on the Web. To do so, we in particular want to detect if a Web page contains completeness statements in natural language, and transform them into RDF-based completeness statements. Last, we also want to increase the expressivity of queries, say, to also handle negations. Queries with negations are especially interesting since negation naturally needs complete information to work correctly.

Acknowledgments

We would like to thank Sebastian Rudolph for his feedback on an earlier version of this paper. The research was supported by the projects “CANDy: Completeness-Aware Querying and Navigation on the Web of Data” and “TaDaQua - Tangible Data Quality with Object Signatures” of the Free University of Bozen-Bolzano, and “MAGIC: Managing Completeness of Data” of the province of Bozen-Bolzano.

References

1. Patrick J. Hayes and Peter F. Patel-Schneider, editors. *RDF 1.1 Semantics*. W3C Recommendation, 25 February 2014.

2. Denny Vrandečić and Markus Krötzsch. Wikidata: A Free Collaborative Knowledgebase. *Commun. ACM*, 57(10):78–85, 2014.
3. Fariz Darari, Werner Nutt, Giuseppe Pirrò, and Simon Razniewski. Completeness Statements about RDF Data Sources and Their Use for Query Answering. In *ISWC 2013*, pages 66–83. Springer Berlin Heidelberg, 2013.
4. Simon Razniewski, Flip Korn, Werner Nutt, and Divesh Srivastava. Identifying the Extent of Completeness of Query Answers over Partially Complete Databases. In *ACM SIGMOD 2015*, pages 561–576, 2015.
5. Steve Harris and Andy Seaborne, editors. *SPARQL 1.1 Query Language*. W3C Recommendation, 21 March 2013.
6. Richard Y. Wang and Diane M. Strong. Beyond Accuracy: What Data Quality Means to Data Consumers. *J. of Management Information Systems*, 12(4):5–33, 1996.
7. Amihai Motro. Integrity = Validity + Completeness. *ACM Trans. Database Syst.*, 14(4), 1989.
8. Alon Y. Levy. Obtaining Complete Answers from Incomplete Databases. In *VLDB 1996*, pages 402–412, 1996.
9. Simon Razniewski and Werner Nutt. Completeness of Queries over Incomplete Databases. *PVLDB*, 4(11):749–760, 2011.
10. Simon Razniewski and Werner Nutt. Assessing Query Completeness over Incomplete Databases. In *VLDB Journal (submitted)*.
11. Christian Fürber and Martin Hepp. SWIQA - a Semantic Web Information Quality Assessment Framework. In *ECIS 2011*, 2011.
12. Pablo N. Mendes, Hannes Mühleisen, and Christian Bizer. Sieve: Linked Data Quality Assessment and Fusion. In *EDBT/ICDT Workshops*, pages 116–123, 2012.
13. Xu Chu, John Morcos, Ihab F. Ilyas, Mourad Ouzzani, Paolo Papotti, Nan Tang, and Yin Ye. KATARA: A data cleaning system powered by knowledge bases and crowdsourcing. In *ACM SIGMOD 2015*, pages 1247–1261, 2015.
14. Maribel Acosta, Elena Simperl, Fabian Flöck, and Maria-Esther Vidal. HARE: A hybrid SPARQL engine to enhance query answers via crowdsourcing. In *K-CAP 2015*, pages 11:1–11:8, 2015.
15. Luis Antonio Galárraga, Christina Teflioudi, Katja Hose, and Fabian M. Suchanek. AMIE: Association Rule Mining under Incomplete Evidence in Ontological Knowledge Bases. In *WWW 2013*, pages 413–422, 2013.
16. Xin Dong, Evgeniy Gabrilovich, Jeremy Heitz, Wilko Horn, Ni Lao, Kevin Murphy, Thomas Strohmann, Shaohua Sun, and Wei Zhang. Knowledge vault: a web-scale approach to probabilistic knowledge fusion. In *ACM SIGKDD 2014*, pages 601–610, 2014.
17. Fariz Darari, Radityo Eko Prasojo, and Werner Nutt. Expressing No-Value Information in RDF. In *ISWC Posters and Demos*, 2015.