

A Framework for Efficient Representative Summarization of RDF Graphs

Šejla Čebirić¹, François Goasdoué^{2,1}, and Ioana Manolescu¹

¹ INRIA, France

² Univ. Rennes 1, France

Abstract. RDF is the data model of choice for Semantic Web applications, and it is often used to model large and heterogeneous datasets. We consider the problem of determining as quickly as possible if a query q lacks answers on a given RDF graph G : from a user perspective, this allows disregarding graphs uninteresting for them; from a query processing perspective, this saves system resources as a query which is known to lack answers does not need to be evaluated on G .

To address this, we introduce a generic *RDF summarization framework* whereas from an RDF G , we build a summary $G_{/\equiv}$, also an RDF graph but often many orders of magnitude smaller, and such that *if q lacks answers on G , it is guaranteed to lack answers on $G_{/\equiv}$* . Further, in the presence of RDF Schema, the interesting question is whether q has answers on the *saturation* of G , denoted G^∞ . One could address that by saturating G (if not already done), then computing $(G^\infty)_{/\equiv}$ and evaluating q on it. For more efficiency, we introduce a *shortcut procedure* by which we compute the summary of (G^∞) *directly from G* without saturating it, and provide a sufficient condition for an RDF summary defined according to our framework to admit this shortcut.

1 Summarization framework

This section recalls the RDF summarization framework defined in [?,?] on which we build the present paper.

Definition 1. *RDF node equivalence* Let \equiv be a binary relation between the nodes of an RDF graph. We say \equiv is an RDF node equivalence relation (or RDF equivalence, in short) iff (i) \equiv is an equivalence relation in the classical sense (it is reflexive, symmetric and transitive), (ii) any class node is \equiv only to itself, and (iii) any property node is \equiv only to itself.

An RDF summary is defined as a graph quotient with respect to a given RDF node equivalence:

Definition 2. *RDF summary* Given an RDF graph G and an RDF node equivalence relation \equiv , the summary of G by \equiv , which is an RDF graph denoted $G_{/\equiv}$, is the quotient of G by \equiv . $G_{/\equiv}$ data nodes use fresh URIs to identify the sets of equivalent G data nodes.

Importantly, any RDF summary enjoys the following property:

Proposition 1. *Schema preservation* An RDF graph G and an RDF summary $G_{/\equiv}$ of it have the same schema triples, i.e., $S_G = S_{G_{/\equiv}}$ holds.

For a summary to reflect (*represent*) the input graph, queries having answers on G should also have answers on the summary. Given an RDF query language (dialect) \mathcal{Q} , we define:

Definition 3. *Query-based representativeness* Let G be any RDF graph. $G_{/\equiv}$ is \mathcal{Q} -representative of G if and only if for any query $q \in \mathcal{Q}$ such that $q(G^\infty) \neq \emptyset$, we have $q((G_{/\equiv})^\infty) \neq \emptyset$.

Representativeness is desirable as the summary can be used to help users formulate queries: therefore, it is important to reflect all graph patterns that may occur in the data.

We focus on the following query language:

Definition 4. *RBGP* queries* An extended relational (RBGP*, in short) query is a BGP query whose body has (i) URIs or variables in all the property positions, (ii) a URI in the object position of every τ triple, and (iii) variables in any other positions.

\implies Both languages [SC: "Both" refers to RBGP and RBGP* but RBGP definition is missing here] forbid URIs or literals in subject and object positions, and require that if type triples are specified in the query, the type is known. They differ in that RBGPs require URIs in the property positions, whereas RBGP* also allow variables there. Clearly, RBGPs are a restriction of RBGP*. A sample RBGP* query is:

$$q^*(x_1, x_3) :- x_1 \tau \text{ Book}, x_1 \text{ author } x_2, x_2 y x_3$$

We define *RBGP* representativeness* by instantiating \mathcal{Q} in Definition 3 to RBGP* queries (Definition 4).

Based on the above definitions, we established [?,?]:

Proposition 2. *Summary representativeness* An RDF summary $G_{/\equiv}$ is RBGP*-representative.

Given that the semantics of G is G^∞ , an RBGP* representative summary must reflect both the explicit and the implicit triples of G . A straightforward way to obtain $(G^\infty)_{/\equiv}$ is to compute G^∞ and then summarize it. This is not directly possible when one does not have the right to add triples to G , and when it is possible, it may be time and space-consuming. Further, it has to be maintained when data or schema triples in G change.

To avoid going through the saturation step, we identify a method called *shortcut* (to be defined shortly), which guarantees that we can build RBGP* representative summaries efficiently. It allows constructing an RDF graph *strongly isomorphic* to $(G^\infty)_{/\equiv}$, which turns out to be *also* a summary of G^∞ , as strongly isomorphic graphs are identical up to renaming of their *data* node URIs:

Definition 5. *Strong isomorphism* \simeq A strong isomorphism between two RDF graphs G_1, G_2 , noted $G_1 \simeq G_2$, is an isomorphism which is the identity for the class and property nodes.

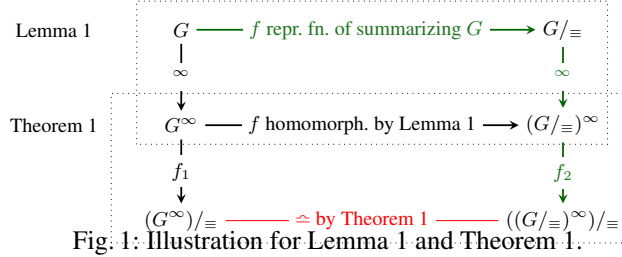


Fig. 1: Illustration for Lemma 1 and Theorem 1.

Definition 6. *Shortcut Summarization through the RDF node equivalence relation \equiv admits a shortcut iff for any RDF graph G , $(G^\infty)_{/\equiv} \simeq ((G/\equiv)^\infty)_{/\equiv}$ holds, where \simeq denotes a strong isomorphism (Definition 5).*

The efficient method (or shortcut) to build $(G^\infty)_{/\equiv}$ introduced above is: summarize G ; saturate the result; then summarize it again. The shortcut leads to a graph whose saturation is strongly isomorphic to that of $(G^\infty)_{/\equiv}$. The two may differ in the exact URIs of the summary *data* nodes (depicted as blank circles in this paper’s examples); these are just *representatives* of G data node groups, and their exact URIs do not matter. In contrast, the summary structure is of crucial importance as representativeness relies on it; this structure is preserved by \simeq . Thus, essentially, *a shortcut allows to obtain (an equivalent to) the summary of the saturated G without saturating it.*

The next theorem establishes a sufficient condition on an RDF node equivalence relation for the existence of a shortcut. To be able to state the theorem, we introduce a new concept and a Lemma.

Representation function f By the summary definition, to every node in G corresponds exactly one node in the summary $G_{/\equiv}$. We call *representation function* and denote $f_{/\equiv}$ (or simply f , when this does not cause confusion) the function associating a summary node to each G node; we say $f(n)$ *represents* n in the summary. An important structural property relates G , G^∞ and the function f :

Lemma 1 (Summarization Homomorphism). *Let G be an RDF graph, $G_{/\equiv}$ its summary and f the corresponding representation function from G nodes to $G_{/\equiv}$ nodes. f defines a homomorphism from G^∞ to $(G_{/\equiv})^\infty$.*

Theorem 1 (Existence of shortcuts). *Given an RDF node equivalence relation \equiv , and an RDF graph G , let $G_{/\equiv}$ be its summary and $f_{/\equiv}$ the corresponding representation function from G nodes to $G_{/\equiv}$ nodes.*

If \equiv satisfies: for any RDF graph G and any pair (n_1, n_2) of G nodes, $n_1 \equiv n_2$ in G^∞ iff $f(n_1) \equiv f(n_2)$ in $(G_{/\equiv})^\infty$, then $(G^\infty)_{/\equiv} \simeq ((G_{/\equiv})^\infty)_{/\equiv}$ holds.

Importantly, *not all RDF equivalence relations admit a shortcut*, as we will illustrate in Section ??.

The condition identified in Theorem 1 is sufficient; finding a necessary (and sufficient) condition is currently open.