



HAL
open science

Query-Oriented Summarization of RDF Graphs

Šejla Čebirić, François Goasdoué, Ioana Manolescu

► **To cite this version:**

Šejla Čebirić, François Goasdoué, Ioana Manolescu. Query-Oriented Summarization of RDF Graphs. [Research Report] RR-8920, INRIA Saclay; Université Rennes 1. 2017. hal-01325900v3

HAL Id: hal-01325900

<https://inria.hal.science/hal-01325900v3>

Submitted on 3 Feb 2017 (v3), last revised 4 Jul 2018 (v6)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Query-Oriented Summarization of RDF Graphs

Šejla Čebirić, François Goasdoué, Ioana Manolescu

**RESEARCH
REPORT**

N° 8920

June 2016

Project-Teams CEDAR



Query-Oriented Summarization of RDF Graphs

Šejla Čebirić, François Goasdoué, Ioana Manolescu

Project-Teams CEDAR

Research Report n° 8920 — version 3 — initial version June 2016 — revised
version February 2017 — 45 pages

Abstract: The Resource Description Framework (RDF) is the W3C’s graph data model for Semantic Web applications. RDF graphs are often large and heterogeneous, thus users may have a hard time getting familiar with the structure and semantics of a graph.

We consider the problem of building *automatically, with no user input, compact RDF graph summaries which represent the complete structure and semantics of the graph, and are representative and accurate for a large useful dialect of SPARQL*. We provide the first technique for building out of an RDF graph a summary of its explicit and implicit data (the latter is due to RDF semantic constraints); a summary for which this is possible is termed *complete*. We introduce four novel summaries and show that two of them are complete. We provide a sufficient condition for RDF summarization completeness, and show that bisimulation-based summaries previously studied satisfy this condition.

We implemented a summarization tool and demonstrate its effectiveness through a set of experiments.

Key-words: Semantic Web, RDF, data summary, inference, reasoning, data compression

RESEARCH CENTRE
SACLAY – ÎLE-DE-FRANCE

1 rue Honoré d'Estienne d'Orves
Bâtiment Alan Turing
Campus de l'École Polytechnique
91120 Palaiseau

Résumés orientés requêtes de graphes RDF

Résumé : RDF est le modèle de données du W3C, fondé sur les graphes, pour les applications du Web Sémantique. Les graphes RDF sont souvent larges et hétérogènes, compliquant la tâche des utilisateurs qui tentent de se familiariser avec leurs structure et leur sémantique. Nous étudions la construction *automatique, sans aucune information fournie par l'utilisateur, de résumés compacts de graphes RDF, qui représente leur structure et sémantique complète, et soient représentatifs et précis pour un dialecte large de requêtes SPARQL*. Nous sommes les premiers à présenter une technique pour construire, à partir d'un graphe RDF, un résumé de ses données explicites et implicites; ces dernières peuvent être présentes dans le graphe en vertu de ses contraintes sémantiques. Nous appelons un résumé pour lequel une telle technique existe, un résumé *complet*. Nous introduisons quatre nouveaux résumés de graphes, et montrons que deux d'entr'eux sont complets. Nous présentons une condition suffisante pour qu'un résumé RDF soit complet, et montrons que les résumés basés sur la bissimulation bidirectionnelle de graphes satisfont cette condition.

Nous avons développé un outil de résumé de graphes, et nous démontrons son intérêt et son efficacité de façon expérimentale, sur des données synthétiques ainsi que sur des données réelles.

Mots-clés : Web Sémantique, RDF, résumé de données, inférence, raisonnement, compression de données

1 Introduction

The Resource Description Framework (RDF) is a graph-based data model promoted by the W3C as the standard for Semantic Web applications. Its associated query language is SPARQL. RDF graphs are *varied*, coming from scientific applications, social or online media, government data etc. They are often *large*, and *heterogeneous*, i.e., resources described in an RDF graph may have very different sets of properties. An RDF resource may have one or several *types* (which may or may not be related to each other) or lack types. *RDF Schema* (RDFS) statements are sometimes part of an RDF graph, in which case they lead to **implicit data**. For instance, if the schema states that anyone *having a driver license* is a *Human* and we know *Bill has the driver license 123*, then the fact that *Bill* is a *Human* is implicitly present in the graph. According to the W3C RDF and SPARQL specification, **the semantics of an RDF graph G , typically denoted G^∞ , comprises both its explicit and implicit data**; in particular, SPARQL query answers must be computed *reflecting all the explicit and implicit data*. These features make an RDF graph complex, both structurally and conceptually; it is intrinsically hard to get familiar with it, especially if a schema is absent or describes only part of the data.

We study the problem of *RDF summarization*, that is: given any input RDF graph G , *efficiently* build a *compact* RDF graph $G_{/\equiv}$ which *summarizes both the structure and semantics of G* . Our summaries can be used as a **GUI** with the help of a graph visualization library, e.g., Dot or Prefuse, and/or to **inform RDF processing tools about the complete graph structure and semantics**; for instance, they can be used to suggest queries to users exploring an unknown RDF graph. Toward this goal, we summarize *without any user input or parameter setting*. The literature features many summaries, e.g., [11, 9, 7] and some focus on RDF, e.g., [34, 21, 27] (see also Section 9). The main novelty of our work is to build from G a summary of **both its explicit and implicit data** (the latter may not be present in G although it holds there). We call such summaries *complete*, and we are the first to define a framework for complete summarization of semantic-rich RDF graphs. Our summaries can be built efficiently, are 2 to 6 orders of magnitude smaller than G , and never larger than G .

The contributions made in this work are as follows:

1. We define *RDF node equivalence relations* taking into account the particular features of RDF, and based on them define *RDF graph summaries as quotient graphs*. We show that any such summary is representative (has some answer for any query that has answers on G) and accurate (any query having answers on the summary has answer on *some* graph which it summarizes¹).
2. We provide a *sufficient condition* on an RDF node equivalence relation, for the RDF summary defined based on this relation to be *complete*. Moreover, we are the first to provide a concrete *method* for summarizing the explicit and implicit data without materializing the implicit data. As our experiments confirm, completeness is desirable as one can *get the full graph summary without spending the time to materialize and the space to store the implicit data*. We adapt the classic bisimulation-based summary definitions from the literature to our framework, and show that the *RDF bisimulation summaries defined in our framework are complete*. This novel result nicely complements existing knowledge about bisimulation summaries e.g., [25, 28, 19, 30, 7].
3. We introduce *four novel RDF summaries*, based on new RDF node equivalence relations (Section 4 and 5). We show that *two of them are complete* (based on our sufficient condition), while the *two others* are not. Interestingly, incompleteness turns out to be due to *separate treatment of RDF nodes having types*; this contrasts with the prominent role played by resource types in previous RDF summarization proposals [34, 21, 27]. We compare our summaries and analyze relationships between them in Section 6.
4. We have devised *summary building algorithms* (Section 7) and implemented our proposal in a summary building tool available for download at [1]. We demonstrate through *experiments* the strong size reduction through summarization for many real and synthetic datasets and the scalability of our summary building algorithms in Section 8.

¹As commonly the case, many graphs may have the same summary.

| Assertion | Triple | Relational notation |
|-----------|-------------------------|---------------------|
| Class | $s \text{ rdf:type } o$ | $o(s)$ |
| Property | $s \text{ p } o$ | $p(s, o)$ |

| Constraint | Triple | OWA interpretation |
|---------------|-------------------------|--------------------------------------|
| Subclass | $s \prec_{sc} o$ | $s \subseteq o$ |
| Subproperty | $s \prec_{sp} o$ | $s \subseteq o$ |
| Domain typing | $p \leftrightarrow_d o$ | $\Pi_{\text{domain}}(s) \subseteq o$ |
| Range typing | $p \leftrightarrow_r o$ | $\Pi_{\text{range}}(s) \subseteq o$ |

Figure 1: RDF (top) & RDFS (bottom) statements.

As an image is worth a thousand words, the Web site [1] provides graphical representations of many sample summaries.

2 Preliminaries

We introduce RDF graphs and queries in Section 2.1, and requirements for our RDF summaries in Section 2.2.

2.1 RDF Graphs and Queries

An *RDF graph* (or *graph*, in short) is a set of *triples* of the form $s \text{ p } o$. A triple states that its *subject* s has the *property* p , and the value of that property is the *object* o . We consider only well-formed triples, as per the RDF specification [37], using uniform resource identifiers (URIs), typed or untyped literals (constants) and blank nodes (unknown URIs or literals). Blank nodes are essential features of RDF allowing to support *unknown URI/literal tokens*. These are conceptually similar to the labeled nulls or variables used in incomplete relational databases [2], as shown in [10].

Notations. We use s , p , and o as placeholders for subjects, properties and objects, respectively. Literals are shown as strings between quotes, e.g., “*string*”. Figure 1 (top) shows how to use triples to describe resources, that is, to express class (unary relation) and property (binary relation) assertions. The RDF standard [37] has a set of *built-in classes and properties*, as part of the `rdf:` and `rdfs:` pre-defined namespaces. We use these namespaces exactly for these classes and properties, e.g., `rdf:type` specifies the class(es) to which a resource belongs. *For brevity, we will sometimes use τ to denote `rdf:type`.*

For example, the following RDF graph G describes a book, identified by `doi1`: its author (a blank node `_:b1` related to the author name), title and date of publication:

$$G = \{ \text{doi}_1 \text{ rdf:type } \text{Book}, \text{ doi}_1 \text{ writtenBy } _ :b_1, \\ \text{doi}_1 \text{ hasTitle } \text{“Le Port des Brumes”}, \\ _ :b_1 \text{ hasName } \text{“G. Simenon”}, \text{ doi}_1 \text{ publishedIn } \text{“1932”} \}$$

RDF Schema (RDFS). RDFS allows enhancing the descriptions in RDF graphs by declaring *semantic constraints* between the classes and the properties they use. Figure 1 (bottom) shows the four main kinds of RDFS constraints, and how to express them through triples. For concision, we denote the properties expressing *subclass*, *subproperty*, *domain* and *range* constraints by the symbols \prec_{sc} , \prec_{sp} , \leftrightarrow_d and \leftrightarrow_r , respectively. Here, “domain” denotes the first, and “range” the second attribute of every property. The RDFS constraints (Figure 1) are interpreted under the *open-world assumption (OWA)* [2], i.e., as *deductive* constraints: if the triples `hasFriend rdfs:domain Person` and `Anne hasFriend Marie` hold in the graph, then so does the triple `Anne rdf:type Person`. The latter is due to the domain constraint in Figure 1.

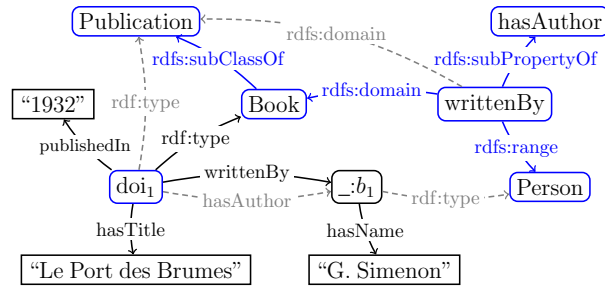


Figure 2: RDF graph and its saturation.

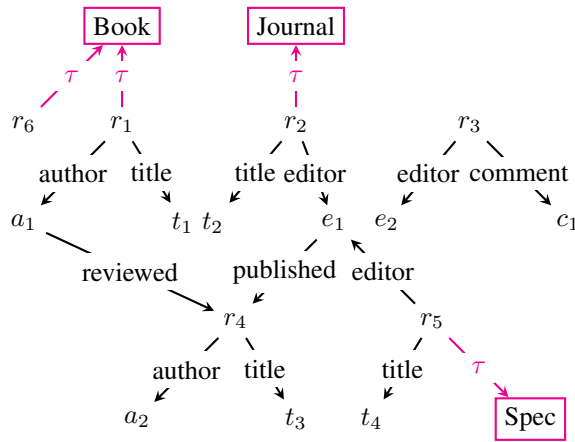


Figure 3: RDF graph used for summarization.

Implicit triples are an important RDF feature, considered part of the RDF graph even though they are not explicitly present in it, e.g., $\text{Anne rdf:type Person}$ above. W3C names *RDF entailment* the mechanism through which, based on a set of explicit triples and some *entailment rules*, implicit RDF triples are derived. We denote by \vdash_{RDF}^i *immediate entailment*, i.e., the process of deriving new triples through a *single* application of an entailment rule. More generally, a triple $s p o$ is entailed by a graph G , denoted $G \vdash_{\text{RDF}} s p o$, if and only if there is a sequence of applications of immediate entailment rules that leads from G to $s p o$ (where at each step, triples previously entailed are also taken into account).

For instance, assume that the RDF graph G above is extended with the constraints: $\text{Book} \prec_{sc} \text{Publication}$, $\text{writtenBy} \prec_{sp} \text{hasAuthor}$, $\text{writtenBy} \leftrightarrow_d \text{Book}$ and $\text{writtenBy} \leftrightarrow_r \text{Person}$.

The resulting graph is depicted in Figure 2. Its implicit triples are those represented by dashed-line edges.

RDF graph saturation. The immediate entailment rules allow defining the finite *saturation* (a.k.a. closure) of an RDF graph G , which is the RDF graph G^∞ defined as the fixpoint obtained by repeatedly applying \vdash_{RDF}^i rules on G .

The saturation of an RDF graph is unique (up to blank node renaming), and does not contain implicit triples (they have all been made explicit by saturation). An obvious connection holds between the triples entailed by a graph G and its saturation: $G \vdash_{\text{RDF}} s p o$ if and only if $s p o \in G^\infty$.

The semantics of an RDF graph is its saturation. RDF entailment is part of the RDF standard itself; in particular, *the answers to a query posed on G must take into account all triples in G^∞* [37].

For presentation purposes, we may use a *triple-based* or a *graph-based* representation of an RDF

graph:

1. The triple-based representation of an RDF graph. We see an RDF graph G as a *partition* of its triples into three components $G = \langle D_G, S_G, T_G \rangle$, where: (i) S_G , the **schema** component, is the set of all G triples whose properties are \prec_{sc} , \prec_{sp} , \leftrightarrow_d or \leftrightarrow_r ; (ii) T_G , the **type** component, is the set of τ triples from G ; (iii) D_G , the **data** component, holds all the remaining triples of G . Note that each of D_G , S_G , and T_G is an RDF graph by itself.

Further, we call *data property* any property p occurring in D_G , and *data triple* any triple in D_G . For illustration, we use the graph in Figure 3: class nodes are shown in purple boxes, while T_G triples appear as purple arrows; D_G consists of the triples shown in black, while S_G is empty. For brevity, we will denote the data properties appearing in Figure 3 by their starting letters: a , t , e , c , r and p .

2. The graph-based representation of an RDF graph. As per the RDF specification [37], the set of *nodes* of an RDF graph is the set of subjects and objects of triples in the graph, while its *edges* correspond to its triples. We define three categories of RDF graph nodes: (i) a *class node* is any node whose URI appears as subject or object of a \prec_{sc} triple, or object of a \leftrightarrow_d or \leftrightarrow_r or τ triple; (ii) a *property node* is any node whose URI appears as subject/object of a \prec_{sp} triple, or subject of a \leftrightarrow_d or \leftrightarrow_r triple, or property of a triple² (iii) a *data node* as any node that is neither a class nor a property node. Note that the sets of class nodes and of property nodes may intersect (indeed, nothing in the RDF specification forbids it), while class nodes are disjoint from data nodes, and property nodes are disjoint from data nodes, too.

Size and cardinality notations. We denote by $|G|_n$ the number of nodes in a graph G , and by $|G|$ its number of edges. Further, for a given attribute $x \in \{s, p, o\}$ and graph G , we note $|G|_x^0$ the number of distinct values of the attribute x within G . For instance, $|D_G|_p^0$ is the number of distinct properties in the data component of G .

Queries. We consider the SPARQL dialect consisting of *basic graph pattern* (BGP) queries, a.k.a. conjunctive queries, widely considered in research but also in real-world applications [20]. A BGP is a set of *triple patterns*, or triples in short; each triple has a subject, property and object, some of which can be variables.

Notations. We use the usual conjunctive query notation $q(\bar{x}) :- t_1, \dots, t_\alpha$, where $\{t_1, \dots, t_\alpha\}$ are triple patterns; the query head variables \bar{x} are called *distinct variables*, and are a subset of those in t_1, \dots, t_α . For boolean queries, \bar{x} is empty. The head of q is $q(\bar{x})$, its body is t_1, \dots, t_α ; x, y, z , etc. denote variables.

Query answering. The *evaluation* of a query q against G has access only to the explicit triples of G , thus it may fail to return the complete answer; the latter is obtained by evaluating q against G^∞ . For instance, the query below asks for name of the author of “Le Port des Brumes”:

$$q(x_3) :- x_1 \text{ hasAuthor } x_2, x_2 \text{ hasName } x_3 \\ x_1 \text{ hasTitle "Le Port des Brumes"}$$

Its answer against the explicit and implicit triples of our sample graph is: $q(G^\infty) = \{\langle \text{"G. Simenon"} \rangle\}$.

Note that evaluating q only against G leads to the empty answer, which is obviously incomplete. This is why, in keeping with the RDF and SPARQL standard, our work considers query *answering*.

2.2 Principles and Design Choices

We decide that *the summary* G_{\equiv} of an RDF graph G is an RDF graph itself. This gives a natural basis for completeness, as RDF schema statements can be part of the summary; it also allows to use existing RDF tools to manipulate the summary. Further, we require summaries to satisfy the following conditions:

²A property node must be a *node*, i.e., merely appearing in a property position does not make an URI a property node; it also needs to appear a subject or object in the same or another triple.

- *Schema independence* It must be possible to summarize G whether or not it has RDF schema information.
- *No user input, no parameters* Summarization must not depend user input nor on manually-specified parameters, as the summary should help users get familiar with the data (and not assume this is already the case).
- *Semantic completeness* The summary of G must reflect all the explicit data and the implicit data, given that the semantic of G is its saturation.
- *Tolerance to heterogeneity* The summary should enable the recognition of “similar” resources despite some amount of heterogeneity in their properties and/or types

The following properties are of a more quantitative nature:

- *Compactness* The summary should be typically smaller than the RDF graph, ideally by many orders of magnitude.
- *Representativeness* The summary should preserve as much information as possible about the structure G^3 .
- *Accuracy* The summary should avoid, to the extent possible, reflecting data that does not exist in G .

To meet our goals of completeness and representativeness, we will not represent frequent graph structures differently / at the expense of the less frequent ones. Instead, we need a logical form of representation, based on queries, as we explain below.

Criteria for representativeness and accuracy. For *representativeness*, queries having answers on G should also have answers on the summary. Symmetrically, for *accuracy*, a query that can be answered on the summary, should also be answerable on the RDF graph itself.

Given an *RDF query language (dialect)* Q , we define:

Definition 1. (QUERY-BASED REPRESENTATIVENESS) *Let G be any RDF graph. $G_{/\equiv}$ is Q -representative of G if and only if for any query $q \in Q$ 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.

To define accuracy, recall that *several graphs may have the same summary*, as summarization loses some information from the input graph; if two RDF graphs differ only with respect to that information, they have the same summary. We term *inverse set* of a summary $G_{/\equiv}$, the set of all RDF graphs whose summary is $G_{/\equiv}$ (the inverse set is purely conceptual, i.e., we never need to compute it). This leads to the accuracy criterion:

Definition 2. (QUERY-BASED ACCURACY) *Let G be any RDF graph, $G_{/\equiv}$ its summary, and \mathcal{G} the inverse set of $G_{/\equiv}$. The summary $G_{/\equiv}$ is Q -accurate if for any query $q \in Q$ such that $q((G_{/\equiv})^\infty) \neq \emptyset$, there exists $G' \in \mathcal{G}$ such that $q(G'^\infty) \neq \emptyset$.*

The above characterizes the accuracy of a summary with respect to *any graph it may correspond to*.

Which dialect should we define representativeness and accuracy for? To get compact summaries, we decide they should not include subject and object values (whether they are URIs, literals or blank nodes); this follows the observation that there are many orders of magnitude more distinct subject and objects, than there are distinct properties in an RDF graph [38]. However, we chose to preserve types (classes) and data properties, as these are essential features of an RDF graph. These choices rule out representativeness for queries with URIs or literals in subjects or objects positions. Instead, we identify two large and useful classes of BGPs:

³Clearly, trade-offs exist between compactness and representativeness; we discuss them in Section 6.

Definition 3. (RELATIONAL AND RBGP* QUERIES) A relational BGP (*RBGP*, in short) query is a BGP query whose body has: (i) URIs in all the property positions, (ii) a URI in the object position of every τ triple, and (iii) variables in any other positions.

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.

Both languages 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 is:

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

while a similar RBGP* 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* and *RBGP* accuracy* by instantiating \mathcal{Q} in Definition 1 and Definition 2, respectively, to RBGP* queries (Definition 3).

Discussion: preserving joins on literals While an RBGP (or RBGP*) query cannot enforce that the subject, property or object in a triple is a certain literal or URI (in other words, it cannot express selections), it does express joins. In particular, if the same literal appears in several places in the graph, an RBGP*-representative summary must preserve the joins enabled by these multiple occurrences. For instance, consider the graph:

| | |
|------------------------|--------------------------|
| s_1 rdf:type city | s_1 zipCode "91120" |
| s_2 rdf:type company | s_2 employeeNo "91120" |
| s_3 rdf:type person | s_3 livesIn "91120" |

An RBGP* representative summary *must* have some results to the hypothetical query asking “cities and companies such that the zipcode of the first is the number of employees of the second”, because on the above graph, the query does have results (due to s_1 and s_2). One may find this query meaningless, and argue that its non-emptiness on G is an accident. However, the very similar RBGP query “people living in the zipcode of a city” (which has results based on s_1 and s_3) is meaningful. In general, without human user input, it is hard to detect that the summary should preserve the results of one, and not of the other. Thus, in the remainder of the work, we rely on *RBGP* representativeness*, knowing it preserves literal-based joins which may sometimes be seen as accidental. To disable this, it suffices to redefine BGP (thus, RBGP and RBGP*) query semantics to consider that for query embedding purposes, two occurrences of the same literal in different places of the graph are not equal, and all our framework would adapt to this; we do not pursue this option further. (Note that the above discussion only concerns *joins on literals*; in contrast, joins on URI are a fundamental aspect of RDF (*linked*) data, and we consider they must be preserved by a meaningful summary.)

3 RDF summarization

Keeping in mind the requirements previously stated, to ensure sound-foundation summaries while still keeping open a space of options, we define our summaries through the classical notion of *quotient graph* from graph theory, parameterized by a *node equivalence relation*. We recall this below:

Definition 4. (QUOTIENT GRAPH) Let A be a label set, $G = (V, E)$ be a labeled directed graph whose vertices are V , whose edges are $E \subseteq V \times V \times A$, with labels from A . Let $\equiv \subseteq V \times V$ be an equivalence relation over the nodes of V .

The quotient graph of G using \equiv , denoted G/\equiv , is a labeled directed graph having (i) a node n_S for each set S of equivalent V nodes, and (ii) an edge $n_{S_1} \xrightarrow{a} n_{S_2}$ for some label $a \in A$ iff there exist two nodes $n_1 \in S_1$ and $n_2 \in S_2$ such that the edge $n_1 \xrightarrow{a} n_2 \in E$.

Quotient-based graph summaries have been used in the past. The novelty of our work lies in the equivalence notions we define, and which in particular enable summary completeness w.r.t. explicit and implicit data. We will rely on:

Definition 5. (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.

where class and property nodes were introduced in Section 2.1.

Next, we first identify interesting relations between the data properties of an RDF graph (Section 3.1). Building on this, we define two RDF node equivalence relations based on data properties, as well as another based on node types (Section 3.2); we also reframe the well-known notion of graph bisimulation [15] in our RDF equivalence setting. We formalize RDF summaries and state some of the important properties they all enjoy (Section 3.3); our novel RDF equivalence relations will lead to defining different summaries in Section 4 and 5.

3.1 Data property relationships and cliques

We start by considering *relations between data properties* in a graph G . The simplest relationship is co-occurrence, when a resource is the source and/or target of two data properties. We generalize this into:

Definition 6. (PROPERTY RELATIONS AND CLIQUES) Let p_1, p_2 be two data properties in G :

1. $p_1, p_2 \in G$ are source-related iff either: (i) a data node in G is the subject of both p_1 and p_2 , or (ii) G holds a data node r and a data property p_3 such that r is the subject of p_1 and p_3 , with p_3 and p_2 being source-related.
2. $p_1, p_2 \in G$ are target-related iff either: (i) a data node in G is the object of both p_1 and p_2 , or (ii) G holds a data node r and a data property p_3 such that r is the object of p_1 and p_3 , with p_3 and p_2 being target-related.

A maximal set of data properties in G which are pairwise source-related (respectively, target-related) is called a source (respectively, target) property clique.

For illustration, consider the sample graph in Figure 3. Here, properties a and t are source-related due to r_1 (condition (i) in the definition). Similarly, t and e are source-related due to r_2 ; consequently, a and e are source-related (condition (ii)). Properties r and p are target-related due to r_4 . The non-empty source cliques are $SC_1 = \{a, t, e, c\}$, $SC_2 = \{r\}$ and $SC_3 = \{p\}$, whereas the non-empty target cliques are $TC_1 = \{a\}$, $TC_2 = \{t\}$, $TC_3 = \{e\}$, $TC_4 = \{c\}$ and $TC_5 = \{r, p\}$.

A source clique can be seen as “all the data properties of same-kind resources”; above, it makes sense to group together properties corresponding to various kinds of publications, such as author, title, editor, and comment. Similarly, a target clique comprises “the data properties of which same-kind resources are values”.

It is easy to see that the set of non-empty source (or target) property cliques is a partition over the data properties of G . Further, if a resource $r \in G$ has data properties, they are all in the same source clique; similarly, all the properties of which r is a value are in the same target clique. This allows us to refer to *the source (or target) clique of r* , denoted $SC(r)$ and $TC(r)$. If r is not the value of any property (respectively, has no property), we consider the target (respectively, source) clique of r to be \emptyset .

The target and source cliques of the resources in the graph shown in Figure 3 are shown in Table 1.

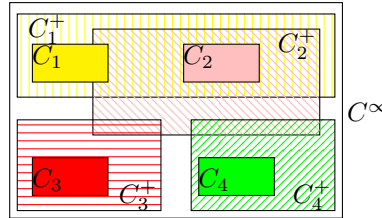
| | | | | | |
|---------|-------------|-------------|-------------|-------------|-------------|
| node | r_1 | r_2 | r_3 | r_4 | r_5 |
| $SC(r)$ | SC_1 | SC_1 | SC_1 | SC_1 | SC_1 |
| $TC(r)$ | \emptyset | \emptyset | \emptyset | TC_5 | \emptyset |
| node | a_1 | t_1 | t_2 | e_1 | e_2 |
| $SC(r)$ | SC_2 | \emptyset | \emptyset | SC_3 | \emptyset |
| $TC(r)$ | TC_1 | TC_2 | TC_2 | TC_3 | TC_3 |
| node | c_1 | t_4 | a_2 | t_3 | r_6 |
| $SC(r)$ | \emptyset | \emptyset | \emptyset | \emptyset | \emptyset |
| $TC(r)$ | TC_4 | TC_2 | TC_1 | TC_2 | \emptyset |

Table 1: Source and target cliques of the sample RDF graph.

Definition 7. (PROPERTY DISTANCE IN A CLIQUE) *The distance between two data properties p, p' in a source (resp. target) clique is:*

- 0 if there exists a resource in G that is the subject (resp. object) of both;
- otherwise, the smallest integer n such that G holds resources $r_0, \dots, r_n \in G$ and data properties p_1, \dots, p_n such that r_0 is the subject (resp. object) of p and p_1, r_1 is the subject (resp. object) of p_1 and p_2, \dots, r_n is the subject (resp. object) of p_n and p' .

In Figure 3, the distance between a and t is 0 since r_1 has both. The distance between a and e is 1, while the distance between a and c is 2. Clearly, p and p' are at distance n for $n > 0$ iff a resource has both p and p'' , and further p'' is at distance $n - 1$ from p' .

Figure 4: Sample cliques of G , their saturations, and their enclosing clique in G^∞ .

Source and target cliques are defined w.r.t. a given graph G . What is the *impact of saturating G on the cliques?* In G^∞ , every G resource has all the data properties it had in G , therefore two data properties belonging to a G clique are also in the same clique of G^∞ . Further, if the schema of G comprises \prec_{sp} constraints, a resource may have in G^∞ a data property that it did not have in G . In turn, this leads to G^∞ cliques which “subsume and fuse” several G cliques into one.

For a given clique C of G , we call *saturated clique* and denote C^+ the of all the properties in C and all their generalizations (superproperties). Observe that C^+ reflects only C and the schema of G ; in particular, it does not reflect the data properties shared by resources in G^∞ , and thus in general C^+ is *not* a clique of G^∞ . The following Lemma characterizes the relationships between cliques of G , their saturated versions C^+ , and cliques of G^∞ :

Lemma 1 (Saturation vs. property cliques). *Let C, C_1, C_2 be distinct non-empty source (or target) cliques of G .*

1. *There exists exactly one source (resp. target) clique C^∞ corresponding to G^∞ such that $C \subseteq C^\infty$.*

2. If $C_1^+ \cap C_2^+ \neq \emptyset$ then all the properties in C_1 and C_2 are in the same G^∞ clique C^∞ .
3. Any non-empty source (or target) clique C^∞ is a union of the form $C_1^+ \cup \dots \cup C_k^+$ for some $k \geq 1$, where each C_i is a non-empty source (resp. target) clique of G . Further, the above saturated cliques are all connected through intersections: for any C_i, C_j where $1 \leq i, j \leq k$ with $i \neq j$, there exist some cliques $D_1 = C_i, \dots, D_n = C_j$ in the set $\{C_1, \dots, C_k\}$ such that:

$$D_1^+ \cap D_2^+ \neq \emptyset, \dots, D_{n-1}^+ \cap D_n^+ \neq \emptyset$$

4. Let p_1, p_2 be two data properties in G , whose source (or target) cliques are C_1 and C_2 . Properties p_1, p_2 are in the same source (resp. target) clique C^∞ corresponding to G^∞ if and only if there exist k non-empty source (resp. target) cliques of G , $k \geq 0$, denoted D_1, \dots, D_k such that:

$$C_1^+ \cap D_1^+ \neq \emptyset, D_1^+ \cap D_2^+ \neq \emptyset, \dots, D_{k-1}^+ \cap D_k^+ \neq \emptyset, D_k^+ \cap C_2^+ \neq \emptyset.$$

Proof. We prove the lemma only for source cliques; the proof for the target cliques is very similar.

1. Any resource $r \in G$ having two data properties also has them in G^∞ ; thus, any data properties in the same source clique in G are also in the same source clique in G^∞ . The unicity of C^∞ is ensured by the fact that the source cliques of G^∞ are by definition disjoint.
2. C_1^+ and C_2^+ intersect on property p iff there exist some $p_1 \in C_1$ and $p_2 \in C_2$ which are specializations of the same p (one, but not both, may also be p itself). Independently, we know that there exist $r_1, r_2 \in G$ such that r_1 has p_1 and r_2 has p_2 ; in G^∞ , r_1 has p_1 and p , thus these two properties are in the same G^∞ clique. Similarly, r_2 has p_2 and p , which ensures that p is also in the same G^∞ source clique.
3. Let $\{p_1, \dots, p_k\}$ be the data properties that appear both in G and in C^∞ ; it follows from the saturation rules and the definition of cliques, that $k > 0$. For $1 \leq i \leq k$, let C_i be the G source clique comprising p_i . Applying lemma point 1., $C_i \subseteq C^\infty$ for each $1 \leq i \leq k$. Further, it is easy to see that $C_i^+ \subseteq C^\infty$, since any property that saturation adds to C_i^+ is also added by saturation to C^∞ . Thus, $\bigcup_{1 \leq i \leq k} C_i^+ \subseteq C^\infty$.

Let us now show that $C^\infty \subseteq \bigcup_{1 \leq i \leq k} C_i^+$. Let $p \in C^\infty$ be a data property, then there exists a resource r having p in G^∞ . Then, in G , r has a property p' which is either p , or is such that $p' \prec_{sp} p$ in G^∞ . Then, in G^∞ , r has both p and p' , which entails that $p' \in C^\infty$. Therefore, p' is a data property occurring both in C^∞ and in G , therefore p' is one of the properties p_i , for some $1 \leq i \leq k$, that is, $p' \in C_i$, and accordingly, $p \in C_i^+$ due to $p' \prec_{sp} p$.

Thus, any data property $p \in C^\infty$ is part of some C_i^+ .

We must still show that the saturated cliques intersect. If $k = 1$ the statement is trivially true. Suppose $k \geq 2$ and the statement is false. Let \mathcal{C} denote the set $\{C_1, \dots, C_m\}$; the cliques in \mathcal{C} are pairwise disjoint by definition. Let $\mathcal{I} \subseteq \mathcal{C}$ be a *maximal* subset of \mathcal{C} cliques such that the saturations of \mathcal{I} cliques all intersect (directly or indirectly). Let $\mathcal{J} = \mathcal{C} \setminus \mathcal{I}$ be the complement of \mathcal{I} ; if the last part of 4. is false, \mathcal{J} is not empty. We denote \mathcal{I}^+ , respectively \mathcal{J}^+ , the set of the saturated cliques from \mathcal{I} , resp. \mathcal{J} .

No data property p_i from \mathcal{I}^+ can be source-related in G^∞ to any data property p_j from \mathcal{J}^+ . This is because source-relatedness requires a resource r having in G^∞ both p_i and a property p source-related to p_j . If such a property p existed, it would belong both to \mathcal{I}^+ (since p has a common source with p_i) and to \mathcal{J}^+ (since p is source-related to p_j); or, \mathcal{I}^+ and \mathcal{J}^+ have no property in common.

The lack of source-relatedness in G^∞ between p_i and p_j chosen as above contradicts the hypothesis that they are part of the same source clique of G^∞ , namely C^∞ .

4. The statement follows quite directly as a consequence of the previous one, concluding our proof.

□

The lemma is crucial as it states the relationship between the cliques (thus, the summaries) of G and G^∞ ; our goal is to build the summary of the latter *without* having access to it, but just access to G . Figure 4 illustrates the lemma. The clique C^∞ of G^∞ encloses the cliques C_1, C_2, C_3, C_4 of G (Lemma item 1). C_2^+ intersects all of C_1^+, C_3^+ and C_4^+ , thus they are all in the same clique C^∞ (item 2), which is the union of C_1^+, C_2^+, C_3^+ and C_4^+ (item 3). A chain of intersecting cliques connects e.g., C_1^+ through C_2^+ to C_4^+ . Item 4 is not illustrated, as individual properties do not show; it follows quite directly from item 3.

3.2 Equivalence relations among graph nodes

We are now ready to introduce two crucial notions of RDF node equivalence based on property cliques:

Definition 8. (STRONG AND WEAK EQUIVALENCE) *Strong (denoted \equiv_s) and weak (denoted \equiv_w) equivalence are RDF node equivalence relations. Two data nodes of G are strongly equivalent, denoted $n_1 \equiv_s n_2$, iff they have the same source and target cliques.*

Two data nodes are weakly equivalent, denoted $n_1 \equiv_w n_2$, iff: (i) they have the same non-empty source or non-empty target clique, or (ii) they both have empty source and empty target cliques, or (iii) they are both weakly equivalent to another node of G .

Observe that strong equivalence implies weak equivalence.

In Figure 3, the resources r_1, r_2, r_3, r_5 are *strongly* equivalent to each other, as well as t_1, t_2, t_3, t_4 . Moreover, r_1, \dots, r_5 are *weakly* equivalent to each other due to their common source clique SC_1 , as well as t_1, t_2, t_3, t_4 due to their common target clique; the same holds for a_1 and a_2 , and separately for e_1 and e_2 . In general, weakly equivalent resources can be connected as exemplified in Figure 5, by an *alternating sequence of source and target cliques*.

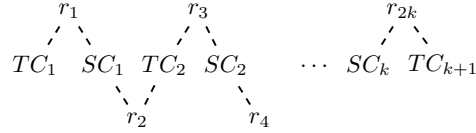


Figure 5: Weakly equivalent resources of G and their cliques.

RDF bisimulation-based node equivalence The classic concept of graph bisimulation has been intensively used to define graph summaries, from [25] to [7] and many others. To clarify the connection with our work, we revisit bisimulation as an RDF equivalence relation, as follows:

Definition 9. (FW RDF BISIMILARITY) *Forward (FW) RDF bisimilarity is an RDF node equivalence relation such that two data nodes $n_1, n_2 \in G$ are FW-bisimilar, denoted $n_1 \equiv_{fw} n_2$, iff (i) for every G data triple of the form $n_1 \text{ p } m_1$ there exists a G data triple of the form $n_2 \text{ p } m_2$ such that $m_1 \equiv_{fw} m_2$, and conversely (ii) for every G data triple of the form $n_2 \text{ p } m_2$ there exists a G data triple of the form $n_1 \text{ p } m_1$ such that $m_1 \equiv_{fw} m_2$.*

Similarly, a backward (BW) RDF bisimilarity, noted \equiv_{bw} , is defined by inverting the subject and object positions of triples in the definition above. Further, a forward and backward (FB) RDF bisimilarity between data nodes, noted \equiv_{fb} is defined by nodes that are both FW and BW similar.

In Figure 3, nodes t_1, t_2, t_4 are BW-bisimilar, as they are only targets of the title property from BW-similar sources r_1, r_2, r_5 , since they have no incoming edges. However, no two resources among r_1 to r_5 are FW-bisimilar (thus, they are not FB-bisimilar either), as they all have different sets of outgoing properties. This example illustrates that *bisimilarity is rare in heterogeneous graphs*, and accordingly, *bisimulation summaries are almost as large at the data e.g., 70% to 99% reported in [19]*. Thus,

bisimulation summaries do not satisfy our compactness requirement. Several works considered *bounded* bisimilarity, i.e., accept nodes as bisimilar if they surrounding subgraphs of *diameter at most k edges* are bisimilar; small values of k reduce summary size, but then (i) structural and semantic properties are only guaranteed on a distance of at most k , (ii) it is hard (and lengthy) for users to guess k through repeated summarization.

The two novel (weak and strong) and three recycled (FW, BW and FB) equivalences mentioned above ignore possible type triples. The last very simple RDF node equivalence we introduce is:

Definition 10. (TYPE-BASED EQUIVALENCE) *Type-based equivalence (denoted \equiv_{τ}) between graph nodes is an RDF equivalence defined as follows. Two data nodes n_1 and n_2 of G are type-equivalent, denoted $n_1 \equiv_{\tau} n_2$, iff they have exactly the same set of types, which is non-empty.*

Recall that some or even all nodes in an RDF graph may lack types; all such nodes are equivalent from the viewpoint of \equiv_{τ} .

Clearly, each equivalence relation defines a partition over the data nodes in G .

3.3 RDF summaries

With the RDF equivalence notions in place, we define RDF summaries directly as graph quotients:

Definition 11. (RDF SUMMARY) *Given an RDF graph G and an RDF relation \equiv among the nodes of G , the summary of G by \equiv , denoted $G_{/\equiv}$, is the quotient of G by the RDF node equivalence relation \equiv .*

Thus, any RDF node equivalence (in particular, any of the six: \equiv_w , \equiv_s , \equiv_{fw} , \equiv_{bw} , \equiv_{fb} and \equiv_{τ} introduced in the previous section) leads to a (possibly distinct) quotient summary. It immediately follows from the definition that *the size of a summary is bounded by that of its summarized graph*. In practice, as we will show, our summaries are much smaller than the graph.

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.

Importantly, any RDF summary enjoys the following property:

Property 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.*

An RDF summary preserves the schema of the summarized graph because (i) class and property nodes, *RDF-equivalent only to themselves*, are preserved through a quotient that uses an RDF node equivalence relation, and (ii) schema triples only contain class or property nodes as subjects/objects. As we shall see, preserving the schema is necessary (but not sufficient) to achieve summary completeness (Section 2.2). Further, since the summary is an RDF graph, the subjects of type and data triples in $G_{/\equiv}$ must be URIs.

Further, we find that:

Proposition 1. (SUMMARY REPRESENTATIVENESS) *An RDF summary $G_{/\equiv}$ is RBGP*-representative.*

Proof. We prove the statement for **RBGP** queries; Proposition 2 (below) carries the statement over to **RBGP***.

Let q be a query such that $q(G^{\infty}) \neq \emptyset$; we need to show that $q((G_{/\equiv})^{\infty}) \neq \emptyset$.

Let $\phi : q \rightarrow G^{\infty}$ be an embedding, assigning to each query variable v , a node from G^{∞} ; we extend ϕ to say it maps triple patterns from q into triples from G^{∞} . We need to produce an embedding from q into $(G_{/\equiv})^{\infty}$.

First, consider a triple pattern t of q whose embedding $\phi(t) \in G^{\infty}$ also belongs to G .

- If $\phi(t)$ is a schema triple, then $\phi(t)$ is also in $G_{/\equiv}$, since G and $G_{/\equiv}$ have the same schema triples.
- Else if $\phi(t)$ is a type triple of the form $s \tau c$, the triple $f(s) \tau c$ belongs to $G_{/\equiv}$, thus also to $(G_{/\equiv})^\infty$.
- Otherwise, $\phi(t) = s p o \in G$ is a data triple, $G_{/\equiv}$ holds the triple $f(s) p f(o)$, thus $(G_{/\equiv})^\infty$ also comprises it.

Now consider a triple pattern t' of q whose embedding $\phi(t') \in G^\infty$ does not belong to G .

- If $\phi(t')$ is a *schema* triple, then by definition of RDF entailment $\phi(t') \in (S_G)^\infty$, thus it is also in $(G_{/\equiv})^\infty$ since $S_G = S_{G_{/\equiv}} \subseteq G_{/\equiv}$ by definition of a summary.
- Else if $\phi(t')$ is a *data* triple in G^∞ , then by definition of RDF entailment, this triple $\phi(t')$ must be entailed by a *data* triple $t_d = s_d p_d o_d$ in G and a subproperty constraints t_s , i.e., a schema triple in S^∞ . As explained above, $f(s_d) p_d f(o_d) \in G_{/\equiv}$; at the same time, t_s also belongs $(G_{/\equiv})^\infty$, since $S_G = S_{G_{/\equiv}}$ hence $(S_G)^\infty = (S_{G_{/\equiv}})^\infty$. It follows that the inference step which entailed $\phi(t')$ from t_d and t_s in G^∞ also applies on $f(s) p_d f(o_d)$ and t_s in $(G_{/\equiv})^\infty$.
- Otherwise, $\phi(t')$ is a τ triple in G^∞ . This may result either:
 - from a D_G triple $t_d = s_d p_d o_d$ and a triple $t_s \in (S_G)^\infty$, if t_s is a \leftrightarrow_d or \leftrightarrow_r triple. This case is very similar to the one above.
 - from a T_G triple of the form $s \tau c_1$ and a schema triple $t_s = c_1 \prec_{sc} c_2 \in (S_G)^\infty$, such that $\phi(t') = s \tau c_2$. In this case, $G_{/\equiv}$ holds the triple $f(s) \tau c_1$ which is also present in $(G_{/\equiv})^\infty$, thus the same inference step applies in $(G_{/\equiv})^\infty$ to produce $f(s) \tau c_2$, since $S_G = S_{G_{/\equiv}}$ hence $(S_G)^\infty = (S_{G_{/\equiv}})^\infty$.

Thus, any q triple mapped by ϕ into a data G^∞ triple (which may or may not explicitly belong to G) is also mapped into a corresponding triple in $(G_{/\equiv})^\infty$.

To conclude this proof, we now need to show that, in addition to the fact that each q triple that has an embedding in G^∞ has also necessarily an embedding in $G_{/\equiv}^\infty$, if q has an embedding in G^∞ , then q has also an embedding in $(G_{/\equiv})^\infty$. This amounts to show that any two q triples t_1 and t_2 that join and that have an embedding in G^∞ also embed in $(G_{/\equiv})^\infty$ (i.e., q joins are preserved).

- If both $\phi(t_1)$ and $\phi(t_2)$ are schema triples in G^∞ , then these two triples are also in $(G_{/\equiv})^\infty$, since $S_G = S_{G_{/\equiv}}$ hence $(S_G)^\infty = S_{G_{/\equiv}}^\infty$.
- Else if both $\phi(t_1)$ and $\phi(t_2)$ are non-schema triples in G^∞ :
 - If both $\phi(t_1)$ and $\phi(t_2)$ are data triples in G^∞ , there exists a triple t'_1 (resp t'_2) in G , with same subject/object, from which $\phi(t_1)$ (resp. $\phi(t_2)$) is entailed using a sub-property constraint t_s^1 (resp. t_s^2) from S^∞ . Since $\phi(t_1)$ and t'_1 (resp. $\phi(t_2)$ and t'_2) have the same subject and object values, then t'_1 and t'_2 have same values on the places where t_1 and t_2 join. Therefore, if we assume that $t'_1 = s_1 p_1 o_1$ and $t'_2 = s_2 p_2 o_2$, the $G_{/\equiv}$ triples $f(s_1) p_1 f(o_1)$ and $f(s_2) p_2 f(o_2)$ necessarily have same values on the places where t_1 and t_2 join. Moreover, since $S_G = S_{G_{/\equiv}}$ hence $(S_G)^\infty = (S_{G_{/\equiv}})^\infty$, these two $G_{/\equiv}$ triples and the above-mentioned t_s^1 and t_s^2 schema triples, produce the counterpart triples of $\phi(t_1)$ and $\phi(t_2)$ in $(G_{/\equiv})^\infty$, which have same subject and object values. Thus, the q triples t_1 and t_2 embed in these two $(G_{/\equiv})^\infty$ triples, if they embed in $\phi(t_1)$ and $\phi(t_2)$ in G^∞ .

- Else if both $\phi(t_1)$ and $\phi(t_2)$ are type triples in G^∞ , say $\phi(t_1) = u \tau c_1$ and $\phi(t_2) = u \tau c_2$, then $t_1 = x \tau c_1$ and $t_2 = x \tau c_2$ by definition of an RBGP query. As in the cases of single q triple embeddings, $\phi(t_1)$ (resp. $\phi(t_2)$) results either from a G triple $u \tau c$ and a S_G^∞ triple $c \prec_{sc} c_1$, or a G triple $u p u_1$ and a S_G^∞ triple $p \leftarrow_d c_1$, or a G triple $u_1 p u$ and a $(S_G)^\infty$ triple $p \hookrightarrow_r c_1$. Therefore, since $S_G = S_{G/\equiv}$ hence $S_G^\infty = (S_{G/\equiv})^\infty$, for $\phi(t_1)$ (resp. $\phi(t_2)$), there are either a $G_{/\equiv}$ triple $f(u) \tau c$ and a $(S_{G/\equiv})^\infty$ triple $c \prec_{sc} c_1$, or a $G_{/\equiv}$ triple $f(u) p f(u_1)$ and a $(S_{G/\equiv})^\infty$ triple $p \leftarrow_d c_1$, or a $G_{/\equiv}$ triple $f(u_1) p f(u)$ and a $(S_{G/\equiv})^\infty$ triple $p \hookrightarrow_r c_1$, which entail $f(u) \tau c_1$ and $f(u) \tau c_2$ in $G_{/\equiv}^\infty$. Thus, the q triples t_1 and t_2 embed in these two $G_{/\equiv}^\infty$ triples, if they embed in $\phi(t_1)$ and $\phi(t_2)$ in G^∞ .
 - Otherwise, $\phi(t_1)$ is a data triple in G^∞ and $\phi(t_2)$ is a type triple in G^∞ . This case is very similar to the two above case, hence we do not detail it.
- Otherwise, $\phi(t_1)$ is a schema triple in G^∞ and $\phi(t_2)$ is not a schema triples in G^∞ . In this case, since q is an RBGP query, q triples must be such that t_1 is $s_1 p_1 o_1$ with $p_1 \in \{\prec_{sc}, \prec_{sp}, \leftarrow_d, \hookrightarrow_r\}$ and t_2 is either $s_2 p o_2$ or $s_2 \tau c$.

Since G^∞ and $G_{/\equiv}^\infty$ have the same schema, $\phi(t_1)$ also belongs to $G_{/\equiv}^\infty$. Now:

- If t_2 is $s_2 p o_2$ then $\phi(t_2) = s, p, o$ must be either in G or entailed from a G data triple s, p', o and a S_G^∞ sub-property triple p', \prec_{sp}, p (see above, for single q triple embedding). If $\phi(t_2)$ is in G , then $f(s) p f(o)$ is in $G_{/\equiv}$, hence in $G_{/\equiv}^\infty$. Otherwise, $f(s) p' f(o)$ is in $G_{/\equiv}$, $p' \prec_{sp} p$ is in $S_{G/\equiv}$ (since G and $G_{/\equiv}$ have the same schema), thus $f(s) p f(o)$ is in $G_{/\equiv}^\infty$. Since t_1 and t_2 joins, s and/or p are class/property nodes. If s (resp. o) is a class/property node, then $f(s) = s$ (resp. $f(o) = o$). Hence, $\phi(t_1) \in G_{/\equiv}^\infty$ joins with $f(s) p f(o) \in G_{/\equiv}^\infty$.
- If t_2 is $s_2 \tau c$ then $\phi(t_2) = s \tau c$ must be either in G or entailed from (i) a G data triple $s p o$ and a S_G^∞ triple $p \leftarrow_d c$, or a G data triple $s_1 p s$ and a S_G^∞ triple $p \hookrightarrow_r c$, or (iii) a G type triple $s \tau c'$ and a S_G^∞ triple $c' \prec_{sc} c$ (see above, for single q triple embedding). If $\phi(t_2)$ is in G , then $f(s) \tau c$ is in $G_{/\equiv}$, hence in $G_{/\equiv}^\infty$. Otherwise, $f(s) p f(o)$ or $f(s_1) p f(s)$ or $f(s) \tau c'$ is in $G_{/\equiv}$, and (since G and $G_{/\equiv}$ have the same schema) thus $f(s) \tau c$ is in $G_{/\equiv}^\infty$. Since t_1 and t_2 can only join on s_2 , s is a class or property nodes, hence $f(s) = s$. Therefore, $\phi(t_1) \in G_{/\equiv}^\infty$ joins with $f(s) \tau c \in G_{/\equiv}^\infty$.

□

Proposition 2. *RBGP representativeness entails RBGP* representativeness.*

Proof. Let q^* be an RBGP* query which is non-empty on G^∞ and $G_{/\equiv}$ be an RBGP-representative summary of G . We show that non-emptiness of q^* on G^∞ entails its non-emptiness on $(G_{/\equiv})^\infty$.

Given that $q^*(G^\infty)$ is non-empty, there exists at least an embedding of q^* into G^∞ ; let q be the query obtained by replacing in q^* , each variable occurring in the property position by the concrete property matching it in G^∞ . Clearly, q has results on G^∞ , and since $G_{/\equiv}$ is RBGP representative, q also has results on $(G_{/\equiv})^\infty$. Therefore, $q(G_{/\equiv}^\infty) \neq \emptyset$, and given that $q \subseteq q^*$ (query containment), it follows that $q^*((G_{/\equiv})^\infty) \neq \emptyset$. □

One can prove in a very similar fashion that summaries are also representative with respect to *regular path expression queries* such as the property paths introduced in SPARQL 1.1: any regular path expression query whose results are non empty on an RDF graph G also has non-empty results on a summary of G . As a consequence, our summaries also preserve reachability in G : if a path leads from n_1 to n_2 in G , the same path also leads from $f(n_1)$, representing n_1 in the summary, to $f(n_2)$.

We are interested in summaries having the following property:

Definition 12. (FIXPOINT PROPERTY) *A summary $G_{/\equiv}$ has the fixpoint property iff for any graph G , $(G_{/\equiv})_{/\equiv} = G_{/\equiv}$ holds.*

Intuitively, the fixpoint property expresses the fact that a summary cannot be summarized further, i.e., $G_{/\equiv}$ is its own summary. This is desirable as we wish our summaries to be as compact as possible. It turns out that any summary has this property (which easily follows from the fact that the summary is a quotient).

Proposition 3. (SUMMARY FIXPOINT) *An RDF summary based on the $\equiv_s, \equiv_w, \equiv_T, \equiv_{fw}, \equiv_{bw}$ or \equiv_{fb} RDF node equivalence relation has the fixpoint property.*

Proof. We start by some remarks which do not depend on the specific equivalence relation used. Suppose that a summary does not have the fixpoint property, i.e., $G_{/\equiv} \neq (G_{/\equiv})_{/\equiv}$. By the definition of a quotient graph, if $(G_{/\equiv})_{/\equiv}$ and $G_{/\equiv}$ differ, then $(G_{/\equiv})_{/\equiv}$ must have *less nodes* than $G_{/\equiv}$: on one hand, a quotient summary always has at most as many nodes as the summarized graph; on the other hand, if G and $G_{/\equiv}$ had the same number of nodes, they would also have the same number of edges, since the summary edges are completely determined by the summary nodes and the edges of the original graph.

We now split the discussion for our three RDF equivalence relations and the three ones based on bisimulation:

Case \equiv_s : if two $G_{/\equiv}$ nodes n_1, n_2 are \equiv_s -equivalent, then both nodes have the same source and target cliques. In this case, n_1 (resp. n_2) represents G nodes, each of which has the same source and target cliques as n_1 and n_2 . Therefore, by definition of a graph quotient, all G nodes represented by n_1 and n_2 would end up in a single $G_{/\equiv}$ node, a contradiction.

Case \equiv_w : if $(G_{/\equiv})_{/\equiv}$ has less nodes than $G_{/\equiv}$, then there exist two $G_{/\equiv}$ nodes n_1, n_2 that have the same non-empty source or target cliques. In this case, there exists some property common to these same cliques, hence all the G nodes represented by n_1 and n_2 would end up in a single $G_{/\equiv}$ node by the definition of a graph quotient using \equiv_w , a contradiction.

Case \equiv_T : if two $G_{/\equiv}$ nodes n_1, n_2 are \equiv_T -equivalent, then both nodes have the same non-empty type sets. In this case, n_1 (resp. n_2) represents G nodes, every of which has the same non-empty type sets as n_1 and n_2 . Therefore, by the definition of a graph quotient, all G nodes represented by n_1 and n_2 end up in a single $G_{/\equiv}$ node, a contradiction.

Case \equiv_{fw} : if two $G_{/\equiv}$ nodes n_1, n_2 are \equiv_{fw} -equivalent, then for every $G_{/\equiv}$ node m_1 such that $n_1 \text{ p } m_1 \in G_{/\equiv}$ there exists a $G_{/\equiv}$ node m_2 such that $n_2 \text{ p } m_2 \in G_{/\equiv}$ and $m_1 \equiv_{fw} m_2$ and similarly for every $G_{/\equiv}$ node m_2 such that $n_2 \text{ p } m_2 \in G_{/\equiv}$ there exists a $G_{/\equiv}$ node m_1 such that $n_1 \text{ p } m_1 \in G_{/\equiv}$ and $m_1 \equiv_{fw} m_2$. By definition of \equiv_{fw} , $n_1 \text{ p } m_1 \in G_{/\equiv}$ iff $r_i^1 \text{ p } r_i^{1'} \in G$ with n_1 representing r_i^1 and m_1 representing $r_i^{1'}$, for $1 \leq i \leq k$. Similarly, $n_2 \text{ p } m_2 \in G_{/\equiv}$ iff $r_j^2 \text{ p } r_j^{2'} \in G$ with n_2 representing r_j^2 and m_2 representing $r_j^{2'}$, for $1 \leq j \leq l$. It therefore follows that the G nodes r_i^1 and r_j^2 are \equiv_{fw} -equivalent, for $1 \leq i \leq k$ and $1 \leq j \leq l$, hence are represented by a single $G_{/\equiv}$ node n (instead of two nodes n_1 and n_2), a contradiction.

Case \equiv_{bw} : the proof derives directly from the above one by considering $G_{/\equiv}$ triples in which n_1, n_2 are objects (instead of subjects).

Case \equiv_{fb} : the proof derives directly from the two above ones by considering that n_1, n_2 are subjects of $G_{/\equiv}$ triples and also objects of some other $G_{/\equiv}$ triples. \square

This has an immediate interesting consequence, considering that an RDF summary belongs to its own inverse set:

Proposition 4. (ACCURACY) *Any RDF summary with the fixpoint property is accurate.*

Proof. The proof follows from the fact that any summary has a fixpoint property: a summary $G_{/\equiv}$ which is its own summary, corresponds at least to the RDF graph $G_{/\equiv}$ itself. \square

An important structural property relates G , G^∞ and the representation function f :

Proposition 5. (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 an homomorphism from G^∞ to $(G_{/\equiv})^\infty$.*

Proof. We first show that an homomorphism can be established from the node sets of G^∞ to that of $(G_{/\equiv})^\infty$.

Recall from Section 2.1 that RDF saturation with RDFS constraints only adds edges between graph nodes, but does not add nodes. Thus, a node n is in G^∞ iff n is in G . Further, by the definition of our quotient-based summaries (Definition 11), n is in G iff $f(n)$ is in $G_{/\equiv}$. Finally, again by the definition of saturation, $f(n)$ is in $G_{/\equiv}$ iff $f(n)$ is in $(G_{/\equiv})^\infty$.

Therefore, every G^∞ node n maps the $f(n)$ $(G_{/\equiv})^\infty$ node (*).

Next, we show that there is a one-to-one mapping between G^∞ edges and those of $(G_{/\equiv})^\infty$.

If $n_1 p n_2$ is an edge in G^∞ , at least one of the following two situations holds:

- $n_1 p n_2$ is an edge in G . This holds iff $f(n_1) p f(n_2)$ is an edge in $G_{/\equiv}$, by definition of an RDF summary. Finally, if $f(n_1) p f(n_2)$ is an edge in $G_{/\equiv}$, then $f(n_1) p f(n_2)$ is also an edge in $(G_{/\equiv})^\infty$.
- $n_1 p' n_2$ is an edge in G , and $p' \prec_{sp} p$ is in S_G^∞ , thus $n_1 p n_2$ is produced by saturation in G^∞ . In this case, we show similarly to the preceding item that $f(n_1) p' f(n_2)$ is an edge in $(G_{/\equiv})^\infty$, hence $f(n_1) p f(n_2)$ is also an edge added to $(G_{/\equiv})^\infty$ by saturation, since $(G_{/\equiv})^\infty$ and G^∞ have the same (saturated) schema triples (Property 1).

If $n_1 \tau c$ is an edge in G^∞ , at least one of the following two situations holds:

- $n_1 \tau c$ is an edge in G . This holds iff $f(n_1) \tau c$ is an edge in $G_{/\equiv}$, by definition of an RDF summary (recall that $f(c) = c$ for classes). Finally, if $f(n_1) \tau c$ is an edge in $G_{/\equiv}$, then $f(n_1) \tau c$ is also an edge in $(G_{/\equiv})^\infty$.
- $n_1 p n_2$ is an edge in G and $p \leftrightarrow_d c$ (or $p \leftrightarrow_r c$) is in S_G^∞ , thus $n_1 \tau c$ is produced by saturation in G^∞ . In this case, we show similarly as above that $f(n_1) p f(n_2)$ is an edge in $(G_{/\equiv})^\infty$, hence $f(n_1) \tau c$ is also an edge added to $(G_{/\equiv})^\infty$ by saturation, since $(G_{/\equiv})^\infty$ and G^∞ have the same (saturated) schema triples (Property 1).

Therefore, every G^∞ edge $n_1 p n_2$ (resp. $n_1 \tau c$) maps into the $(G_{/\equiv})^\infty$ edge $f(n_1) p f(n_2)$ (resp. $f(n_1) \tau c$) (**).

From (*) and (**), it follows that f is an homomorphism from G^∞ to $(G_{/\equiv})^\infty$. \square

We now consider the completeness requirement introduced in Section 2.2. We want the summary of G to reflect both the explicit and the implicit triples of G , that is, we would like to compute $(G^\infty)_{/\equiv}$ starting from $G_{/\equiv}$ and without saturating G . Recall from Section 2 that the semantics of an RDF graph is its saturation; thus, any graph having *the same saturation as* $(G^\infty)_{/\equiv}$ can be seen as the complete summary of G . Finally, observe that RDF summary nodes are just representatives of the original G nodes, thus the identity of summary nodes is not as important as the summary structure. Therefore, we are interested in computing from G and without saturating G , an RDF graph *whose saturation is isomorphic to the saturation of* $(G^\infty)_{/\equiv}$. We formalize this by:

Definition 13. (SUMMARY COMPLETENESS) *The RDF summary defined by a given RDF node equivalence relation \equiv is complete iff $(G^\infty)_{/\equiv} \sim ((G_{/\equiv})^\infty)_{/\equiv}$, where \sim denotes isomorphism between the saturations of two RDF graphs.*

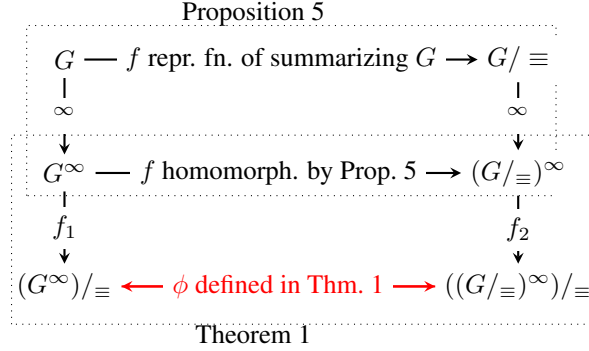


Figure 6: Illustration for Proposition 5 and Theorem 1.

In the above, we used \sim to denote an isomorphism between two saturated RDF graphs, which moreover is the identity on class and property nodes (introduced in Section 2).

Definition 13 gives *the method we found* to compute a graph equivalent to (having the same saturation as) $(G^\infty)_{/\equiv}$ as follows: summarize G ; saturate the result; then summarize it again. In general, we cannot rule out the existence of other methods to find $(G^\infty)_{/\equiv}$ without saturating G . It is, however, worth noting that computing $(G_{/\equiv})^\infty$ (two processing steps starting from G , instead of three) does *not* suffice, for the RDF equivalence notions (and thus, the summaries) that we introduce. This will be exemplified e.g., in Figure 9.

The first major result of our work is:

Theorem 1 (Condition for summary completeness). *Let G be an RDF graph, $G_{/\equiv}$ its summary and f_{\equiv} the corresponding representation function from G nodes to $G_{/\equiv}$ nodes.*

If the RDF node equivalence relation \equiv satisfies: for every 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 the summary of G through \equiv is complete.

Proof. We start by introducing some **notations** (see Figure 6). Let f_1 be the representation function from G^∞ into $(G^\infty)_{/\equiv}$, and f_2 be the representation function from $(G_{/\equiv})^\infty$ into $((G_{/\equiv})^\infty)_{/\equiv}$.

Let the function φ be a function from the $(G^\infty)_{/\equiv}$ nodes to the $((G_{/\equiv})^\infty)_{/\equiv}$ nodes defined as: $\varphi(f_1(n)) = f_2(f(n))$ for n any G^∞ node.

Suppose that for every pair (n_1, n_2) of G nodes, $n_1 \equiv n_2$ in G^∞ iff $f(n_1) \equiv f(n_2)$ in $(G_{/\equiv})^\infty$ holds. Let us show that this condition suffices to ensure $(G^\infty)_{/\equiv} \equiv ((G_{/\equiv})^\infty)_{/\equiv}$ holds, i.e., the φ function defines an isomorphism from $(G^\infty)_{/\equiv}$ to $((G_{/\equiv})^\infty)_{/\equiv}$.

First, let us show that φ is a bijection from all the $(G^\infty)_{/\equiv}$ nodes to all the $((G_{/\equiv})^\infty)_{/\equiv}$ nodes. Since for every pair n_1, n_2 of G^∞ nodes, $n_1 \equiv n_2$ iff $f(n_1) \equiv f(n_2)$ in $(G_{/\equiv})^\infty$, it follows that $(G^\infty)_{/\equiv}$ and $((G_{/\equiv})^\infty)_{/\equiv}$ have the same number of nodes (*).

Further, a given node n in $(G^\infty)_{/\equiv}$ represents a set of equivalent nodes n_1, \dots, n_k from G^∞ . By hypothesis, $n_1 \equiv \dots \equiv n_k$ in G^∞ iff $f(n_1) \equiv \dots \equiv f(n_k)$ in $(G_{/\equiv})^\infty$ holds. Hence, every node $n = f_1(n_1) = \dots = f_1(n_k)$ of $(G^\infty)_{/\equiv}$ maps to a distinct node $n' = f_2(f(n_1)) = \dots = f_2(f(n_k))$ in $((G_{/\equiv})^\infty)_{/\equiv}$ (**).

Similarly, a given node n' in $((G_{/\equiv})^\infty)_{/\equiv}$ represents a set of equivalent nodes $n'_1 = f(n_1), \dots, n'_k = f(n_k)$ in $(G_{/\equiv})^\infty$. By hypothesis, $f(n_1) \equiv \dots \equiv f(n_k)$ in $(G_{/\equiv})^\infty$ iff $n_1 \equiv \dots \equiv n_k$ in G^∞ holds. Hence, every node $n' = f_2(f(n_1)) = \dots = f_2(f(n_k))$ in $((G_{/\equiv})^\infty)_{/\equiv}$ maps to a distinct node $n = f_1(n_1) = \dots = f_1(n_k)$ of $(G^\infty)_{/\equiv}$ (***) .

From (*), (**), and (***), it follows that φ is a bijective function from all the $(G^\infty)_{/\equiv}$ nodes to all the $((G_{/\equiv})^\infty)_{/\equiv}$ nodes.

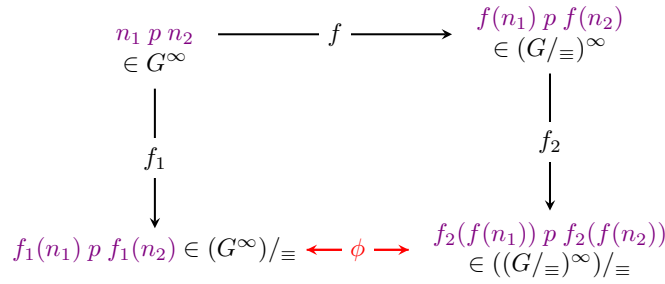


Figure 7: Diagram for the proof of Theorem 1.

Now, let us show that φ defines an isomorphism from $(G^\infty)_{/\equiv}$ to $((G/\equiv)^\infty)_{/\equiv}$.

For every edge $n'_1 p n'_2$ in $(G^\infty)_{/\equiv}$, by definition of an RDF summary, there exists an edge $n_1 p n_2$ in G^∞ such that $n'_1 p n'_2 = f_1(n_1) p f_1(n_2)$. Figure 7 illustrates the discussion. Further, if $n_1 p n_2$ is in G^∞ , then $f(n_1) p f(n_2)$ is in $(G/\equiv)^\infty$ (Proposition 5), hence $f_2(f(n_1)) p f_2(f(n_2))$ is in $((G/\equiv)^\infty)_{/\equiv}$. Therefore,

- since for every $f_1(n_1) p f_1(n_2)$ edge in $(G^\infty)_{/\equiv}$, there is an edge $f_2(f(n_1)) p f_2(f(n_2))$ in $((G/\equiv)^\infty)_{/\equiv}$, and
- since $\varphi(f_1(n)) = f_2(f(n))$, for n any G^∞ node, is a bijective function from all $(G^\infty)_{/\equiv}$ nodes to all $((G/\equiv)^\infty)_{/\equiv}$ nodes,
- it follows that $((G/\equiv)^\infty)_{/\equiv}$ contains the image of all $(G^\infty)_{/\equiv}$ $f_1(n_1) p f_1(n_2)$ triples through φ (*).

Now, for every edge $n''_1 p n''_2$ in $((G/\equiv)^\infty)_{/\equiv}$, by definition of an RDF summary, there exists an edge $n'_1 p n'_2$ in $(G/\equiv)^\infty$ such that $n''_1 p n''_2 = f_2(n'_1) p f_2(n'_2)$. Hence, by Proposition 5, there exists an edge $n_1 p n_2$ in G^∞ such that $n'_1 p n'_2 = f(n_1) p f(n_2)$. Moreover, since $n_1 p n_2$ is in G^∞ , $f_1(n_1) p f_1(n_2)$ is in $(G^\infty)_{/\equiv}$. Therefore, since for every $f_2(f(n_1)) p f_2(f(n_2))$ edge in $((G/\equiv)^\infty)_{/\equiv}$, there is an edge $f_1(n_1) p f_1(n_2)$ in $(G^\infty)_{/\equiv}$, and since $\varphi(f_1(n)) = f_2(f(n))$, for n any G^∞ node, is a bijective function from all $(G^\infty)_{/\equiv}$ nodes to all $((G/\equiv)^\infty)_{/\equiv}$ nodes, $(G^\infty)_{/\equiv}$ contains the image of all $((G/\equiv)^\infty)_{/\equiv}$ $n''_1 p n''_2$ triples through φ^{-1} (**).

Similarly, for every edge $n'_1 \tau c$ in $(G^\infty)_{/\equiv}$, by definition of an RDF summary, there exists an edge $n_1 \tau c$ in G^∞ such that $n'_1 \tau c = f_1(n_1) \tau c$. Further, if $n_1 \tau c$ is in G^∞ , then $f(n_1) \tau c$ is in $(G/\equiv)^\infty$ (Proposition 5), hence $f_2(f(n_1)) \tau c$ is in $((G/\equiv)^\infty)_{/\equiv}$. Therefore,

- since for every $f_1(n_1) \tau c$ edge in $(G^\infty)_{/\equiv}$, there is an edge $f_2(f(n_1)) \tau c$ in $((G/\equiv)^\infty)_{/\equiv}$, and
- since $\varphi(f_1(n)) = f_2(f(n))$, for n any G^∞ node, is a bijective function from all $(G^\infty)_{/\equiv}$ nodes to all $((G/\equiv)^\infty)_{/\equiv}$ nodes,
- it follows that $((G/\equiv)^\infty)_{/\equiv}$ contains the image of all $(G^\infty)_{/\equiv}$ $f_1(n_1) \tau c$ triples through φ (*').

Now, for every edge $n''_1 \tau c$ in $((G/\equiv)^\infty)_{/\equiv}$, by definition of an RDF summary, there exists an edge $n'_1 \tau c$ in $(G/\equiv)^\infty$ such that $n''_1 \tau c = f_2(n'_1) \tau c$. Hence, by Proposition 5, there exists an edge $n_1 \tau c$ in G^∞ such that $n'_1 \tau c = f(n_1) \tau c$. Moreover, since $n_1 \tau c$ is in G^∞ , $f_1(n_1) \tau c$ is in $(G^\infty)_{/\equiv}$. Therefore, since for every $f_2(f(n_1)) \tau c$ edge in $((G/\equiv)^\infty)_{/\equiv}$, there is an edge $f_1(n_1) \tau c$ in $(G^\infty)_{/\equiv}$,

and since $\varphi(f_1(n)) = f_2(f(n))$, for n any G^∞ node, is a bijective function from all $(G^\infty)_{/\equiv}$ nodes to all $((G_{/\equiv})^\infty)_{/\equiv}$ nodes, $(G^\infty)_{/\equiv}$ contains the image of all $((G_{/\equiv})^\infty)_{/\equiv} n_1'' \tau c$ triples through φ^{-1} (**').

From (*) and (**), and, (*') and (**'), it follows that φ defines an isomorphism from $(G^\infty)_{/\equiv}$ to $((G_{/\equiv})^\infty)_{/\equiv}$. \square

The condition is sufficient, and we will use it to establish completeness for our summaries based on strong and weak equivalence; finding a necessary condition is currently open. For the summaries based on bisimulation, we readily establish based on Theorem 1:

Theorem 2. \equiv_{fw} , \equiv_{bw} and \equiv_{fb} summaries are complete.

Proof. We first prove the claim for \equiv_{fw} .

We show this result using the sufficient condition stated in Theorem 1. That is, $n_1 \equiv_{fw} n_2$ in G^∞ holds iff $f(n_1) \equiv_{fw} f(n_2)$ in $(G_{/\equiv_{fw}})^\infty$ holds.

This holds for class nodes and for property nodes since, by definition, they are only equivalent to themselves through some RDF node equivalence relation.

Now, consider two data nodes n_1, n_2 in G^∞ such that $n_1 \equiv_{fw} n_2$ in G^∞ , and let us show that $f(n_1) \equiv_{fw} f(n_2)$ in $(G_{/\equiv_{fw}})^\infty$.

If $n_1 \equiv_{fw} n_2$ holds in G^∞ , then for every triple $n_1 \text{ p } m_1$ there exists a triple $n_2 \text{ p } m_2$ such that $m_1 \equiv_{fw} m_2$ holds, and conversely for every triple $n_2 \text{ p } m_2$ there exists a triple $n_1 \text{ p } m_1$ such that $m_1 \equiv_{fw} m_2$ holds.

Let $\mathcal{P}_{n_1, n_2 \rightarrow m_1, m_2}^\infty$ be the set of outgoing properties from n_1 to m_1 and from n_2 to m_2 in G^∞ .

In G , the set of outgoing properties from n_1 to m_1 , denoted $\mathcal{P}_{n_1 \rightarrow m_1}$ is a subset of $\mathcal{P}_{n_1, n_2 \rightarrow m_1, m_2}^\infty$, since by definition the saturation of a graph only adds edges; similarly, in G , the set of outgoing properties from n_2 to m_2 , denoted $\mathcal{P}_{n_2 \rightarrow m_2}$ is a subset of $\mathcal{P}_{n_1, n_2 \rightarrow m_1, m_2}^\infty$, which may be different from $\mathcal{P}_{n_1 \rightarrow m_1}$.

By definition of a \equiv_{fw} -summary, the set of outgoing properties from $f(n_1)$ to $f(m_1)$ in $G_{/\equiv_{fw}}$ is exactly $\mathcal{P}_{n_1 \rightarrow m_1}$ and similarly the set of outgoing properties from $f(n_2)$ to $f(m_2)$ in $G_{/\equiv_{fw}}$ is exactly $\mathcal{P}_{n_2 \rightarrow m_2}$.

Since G and $G_{/\equiv_{fw}}$ have the same schema (Property 1), it follows that in $(G_{/\equiv_{fw}})^\infty$, the set of outgoing properties from $f(n_1)$ to $f(m_1)$, and from $f(n_2)$ to $f(m_2)$, is exactly $\mathcal{P}_{n_1, n_2 \rightarrow m_1, m_2}^\infty$ (data edges can only be added through \prec_{sp} constraints).

Since the above holds for any pair of data nodes n_1, n_2 such that $n_1 \equiv_{fw} n_2$ in G^∞ , and for any of their G^∞ outgoing edges $n_1 \text{ p } m_1$ and $n_2 \text{ p } m_2$, hence $f(n_1) \equiv_{fw} f(n_2)$ in $(G_{/\equiv_{fw}})^\infty$ holds.

Now, consider two data nodes $f(n_1), f(n_2)$ in $(G_{/\equiv_{fw}})^\infty$ such that $f(n_1) \equiv_{fw} f(n_2)$ in $(G_{/\equiv_{fw}})^\infty$ and let us show that $n_1 \equiv_{fw} n_2$ holds in G^∞ .

If $f(n_1) \equiv_{fw} f(n_2)$ holds in $(G_{/\equiv_{fw}})^\infty$, then for every triple $f(n_1) \text{ p } f(m_1)$ there exists a triple $f(n_2) \text{ p } f(m_2)$ such that $f(m_1) \equiv_{fw} f(m_2)$ holds, and conversely for every triple $f(n_2) \text{ p } f(m_2)$ there exists a triple $f(n_1) \text{ p } f(m_1)$ such that $f(m_1) \equiv_{fw} f(m_2)$ holds.

Let $\mathcal{P}_{f(n_1), f(n_2) \rightarrow f(m_1), f(m_2)}^\infty$ be the set of outgoing properties from $f(n_1)$ to $f(m_1)$ and from $f(n_2)$ to $f(m_2)$ in $(G_{/\equiv_{fw}})^\infty$.

In $G_{/\equiv_{fw}}$, the set of outgoing properties from $f(n_1)$ to $f(m_1)$, denoted $\mathcal{P}_{f(n_1) \rightarrow f(m_1)}$ is a subset of $\mathcal{P}_{f(n_1), f(n_2) \rightarrow f(m_1), f(m_2)}^\infty$, since by definition the saturation of a graph only adds edges; similarly, in $G_{/\equiv_{fw}}$, the set of outgoing properties from $f(n_2)$ to $f(m_2)$, denoted $\mathcal{P}_{f(n_2) \rightarrow f(m_2)}$ is a subset of $\mathcal{P}_{f(n_1), f(n_2) \rightarrow f(m_1), f(m_2)}^\infty$, which may be different from $\mathcal{P}_{f(n_1) \rightarrow f(m_1)}$.

By definition of a \equiv_{fw} -summary, the set of outgoing properties from n_1 to m_1 in G is exactly $\mathcal{P}_{f(n_1) \rightarrow f(m_1)}$ and similarly the set of outgoing properties from n_2 to m_2 in G is exactly $\mathcal{P}_{f(n_2) \rightarrow f(m_2)}$.

Since G and $G_{/\equiv_{fw}}$ have the same schema (Property 1), it follows that in G^∞ , the set of outgoing properties from n_1 to m_1 , and from n_2 to m_2 , is exactly $\mathcal{P}_{f(n_1), f(n_2) \rightarrow f(m_1), f(m_2)}^\infty$ (data edges can only be added through \prec_{sp} constraints).

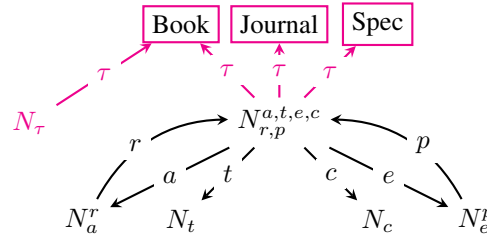


Figure 8: Weak summary of the RDF graph in Figure 3.

Since the above holds for any pair of data nodes $f(n_1), f(n_2)$ such that $f(n_1) \equiv_{f_w} f(n_2)$ in $(G_{/\equiv_{f_w}})^\infty$, and for any of their $(G_{/\equiv_{f_w}})^\infty$ outgoing edges $f(n_1) \text{ p } f(m_1)$ and $f(n_2) \text{ p } f(m_2)$, hence $n_1 \equiv_{f_w} n_2$ in G^∞ holds.

The proof for \equiv_{b_w} directly derives from the above one by considering incoming edges instead of outgoing ones; the proof for \equiv_{t_b} then derives from those of \equiv_{f_w} and \equiv_{b_w} by considering both incoming and outgoing edges. \square

4 Weak-equivalence summaries

In this section, we explore summaries based on the weak equivalence \equiv_w of graph nodes.

4.1 Weak summary

Our first new summary is solely based on weak equivalence:

Definition 14. (WEAK SUMMARY) *The weak summary of the graph G , denoted $G_{/w}$, is its quotient graph w.r.t. the weak equivalence relation \equiv_w .*

For each set of \equiv_w -equivalent G nodes, there is exactly one node in $G_{/w}$. Note that the partition of G nodes into \equiv_w equivalence classes is also a *partition of G data properties at the source*, that is: the sources of G edges labeled with a given data property p are all weakly equivalent. For instance, in Figure 3, the sources of all “editor” edges are in the weak equivalence class $\{r_1, r_2, r_3, r_4, r_5\}$.

The \equiv_w partition over the set of nodes of G also induces a *partition of the source cliques of G* . Indeed, if $n_1 \equiv_w n_2$, their source cliques are connected through a chain of alternating cliques as sketched in Figure 5; conversely, such a chain, by definition, only connects weakly related resources. By a symmetrical reasoning, the weak equivalence classes of G also lead to a *partition over the target cliques of G* . For instance, in Figure 3, to the equivalent resource set $\{r_1, \dots, r_5\}$ corresponds the set of source cliques $\{SC_1\}$ and the set of target cliques $\{TC_5\}$. Thus, *to each set S of weakly equivalent G nodes one can associate through a bijection, a set of G source cliques, and a set of G target cliques*. This leads to:

Definition 15. (SET REPRESENTATIVE) *Let N be any injective function taking as input two sets of URIs, and returning a new URI. Let S be an equivalence class of G data nodes w.r.t. \equiv_w . The weak summary node representing all S nodes, also called the set representative of S in the summary, is: $N(\bigcup_{r \in S} TC(r), \bigcup_{r \in S} SC(r))$.*

In the above, N is called on: the set of data properties from all the target cliques of S nodes; and the set of data properties from all the source cliques of S nodes. We will use N to denote any function which assigns URIs to nodes in quotient (RDF) graphs.

Notations. We use N_r to denote the weak summary node representing a resource $r \in S$, and N_{SC}^{TC} to denote $N(TC, SC)$. For simplicity, we will mostly *omit the set delimiters* when showing TC and SC , and omit one such set altogether if it is empty.

For any resource $r \in G$ which has types and a non-empty source or target clique, its types are carried to N_r in the weak summary.

In the particular case where a data resource $r \in G$ is neither the source nor the target of data properties, i.e., $TC(r) = SC(r) = \emptyset$ (thus r can only appear in τ triples), r is represented by $N(\emptyset, \emptyset)$ which we denote N_τ in the sequel. Observe that if a resource r such that $TC(r) = SC(r) = \emptyset$ has types, the weak summary carries the respective types to N_τ .

The weak summary of the graph in Figure 3 is shown in Figure 8. Its nodes are: $N_{r,p}^{a,t,e,c}$ for the relatedness partition set $\{r_1, \dots, r_5\}$. The target properties of this node are $TC(r_4)$ since the other nodes have empty target clique; the source properties are those in $SC(r_1)$ which is also the source clique of all the other resources in the set; N_a^r for the set $\{a_1, a_2\}$; N_t for the relatedness partition set $\{t_1, t_2, t_3, t_4\}$ etc. The edges from $N_{r,p}^{a,t,e,c}$ to N_a and N_t copy the outgoing edges of r_1 , represented by $N_{r,p}^{a,t,e,c}$; the edge to N_e^p is due to r_2 and r_3 ; the edge to N_c is due to r_3 . The edge from N_a^r to $N_{r,p}^{a,t,e,c}$ is due to a_1 , and the edge from N_e^p to $N_{r,p}^{a,t,e,c}$ is due to e_1 . The τ edges outgoing $N_{r,p}^{a,t,e,c}$ are due to the resources r_1, r_2 and r_3 ; the creation of N_τ (shown in purple font) is due to the node r_6 in the original graph.

The weak summary has the following important properties:

Proposition 6. (UNIQUE DATA PROPERTIES) *Each G data property appears exactly once in $G_{/w}$.*

Proof. First, note that any two weak summary nodes n_1, n_2 cannot be targets of the same data property. Indeed, if such a data property p existed, let TC be the target clique it belongs to. By the definition of the weak summary, n_1 was created by a call to $N(UTC_1, USC_1)$ such that TC is included in the union of target cliques UTC_1 . Similarly, n_2 was created by a call to $N(UTC_2, USC_2)$ such that TC is included in the union of target cliques UTC_2 . Then, $UTC_1 \cap UTC_2 \supseteq TC \neq \emptyset$, which contradicts the fact that different equivalence classes of G nodes correspond to disjoint sets of target cliques. The same holds for the sets of properties of which weak summary nodes are sources. Thus, any data property has at most one source and at most one target in $G_{/w}$. Further, by the definition of the summary as a quotient, every data property present in G also appears in the summary. Thus, there is exactly one p -labeled edge in $G_{/w}$ for every data property in G . \square

Importantly, the above Property 6 warrants that the number of data edges in $G_{/w}$ is exactly $|D_G|_p^0$, the number of distinct data properties in G . Thus, its number of data nodes is at most $2|D_G|_p^0$. The number of type triples in $G_{/w}$ is bound by $\min(|T_G|_e, 2|D_G|_p^0 * |T_G|_o^0)$: the latter corresponds to the case when every data node in $G_{/w}$ is of every type in T_G .

The next important property enjoyed by weak summaries is *completeness*, in the sense of Definition 13. Figure 9 exemplifies this, by tracing the transformation of a graph G on one hand, into $G_{/w}$, $(G_{/w})^\infty$, and then $((G_{/w})^\infty)_{/w}$; and on the other hand, from G to G^∞ then $(G^\infty)_{/w}$. The graph at the bottom right in the figure is at the same time $((G_{/w})^\infty)_{/w}$ and $(G^\infty)_{/w}$.

W -summary completeness is our most technically involved result. To establish it, we rely on the next Proposition:

Proposition 7. (SAME CLIQUES-W) *G^∞ and $(G_{/w})^\infty$ have identical source clique sets, and identical target clique sets. Further, any node n in G^∞ has the source clique SC and target clique TC iff $f(n)$ in $(G_{/w})^\infty$ has the source clique SC and target clique TC , with f the representation function (recall Proposition 5).*

Proposition 7 is established based on Proposition 5, and on the following Lemma:

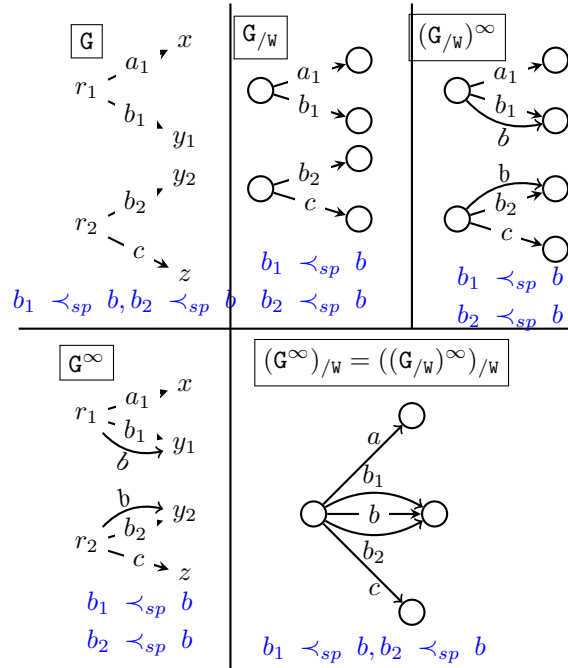


Figure 9: Weak summary completeness illustration (summary nodes shown as unlabeled circles).

Lemma 2. *Data properties are target-related (resp. source-related) in $(G/w)^\infty$ iff they are target-related (resp. source-related) in G^∞ .*

Proof. We prove the lemma for target-related properties.

“Only if”: If data properties are target-related in $(G/w)^\infty$, then they belong to a same target clique TC_W^∞ in $(G/w)^\infty$.

By Lemma 1, point 3, it follows that TC_W^∞ is the union of the saturations of a set of G/w cliques $(TC_W^1)^+, (TC_W^2)^+, \dots, (TC_W^m)^+$. Then:

- For every $1 \leq j \leq m$:
 - TC_W^j is the target clique of a G/w node n^j ;
 - n^j represents a set of weakly-equivalent G resources, which are targets only of properties in TC_W^j . Thus, the properties in TC_W^j are target-related in G .
 - Thus, in G^∞ , also, the properties in TC_W^j are target-related.
 - From this and the definition of a saturated graph and of a saturated target clique, it follows that the properties from $(TC_W^j)^+$ are target-related in G^∞ .
- Further, still by Lemma 1, point 4, each $(TC_W^j)^+$ intersects at least another $(TC_W^l)^+$ for $1 \leq l \neq j < m$, thus the target properties in all the $(TC_W^j)^+$ for $1 \leq j \leq m$, and in particular p , are target-related to each other in G^∞ . Thus, p is target-related in G^∞ to all properties from TC_W^∞ .

“If”: if data properties are target-related in G^∞ , then they belong to a same target clique TC^∞ in G^∞ . Let n_1, \dots, n_k be the set of all G resources which are values of some properties in TC^∞ . By definition of an RDF summary and Proposition 5, each image of $f(n_i) = N(\cup_{r \in S_{n_i}} TC(r), \cup_{r \in S_{n_i}} SC(r))$ of n_i ,

for $1 \leq i \leq k$ is at least the object of the same properties as n_i , hence all the properties of TC^∞ in G^∞ are target-related in $(G/W)^\infty$. \square

The completeness of W -summaries, stated by Theorem 3 below, is then shown based on Proposition 7 by proving that W -summaries enjoy the sufficient condition stated by Theorem 1.

Theorem 3 (W -completeness). *Weak summaries are complete.*

Proof. We show that W summaries enjoy the sufficient condition for completeness stated in Theorem 1. That is, given two nodes n_1, n_2 in G^∞ , f the representation homomorphism corresponding to the weak-equivalence relation \equiv_w , and $f(n_1), f(n_2)$ the images of n_1, n_2 in $(G/W)^\infty$ through f (recall Proposition 5), let us prove that: $n_1 \equiv_w n_2$ in G^∞ iff $f(n_1) \equiv_w f(n_2)$ in $(G/\equiv_w)^\infty$.

Two summary nodes are weakly-equivalent if they have the same non-empty source or target clique. The claim to prove hence immediately follows from Property 7, which guarantees that n_1 and $f(n_1)$, resp. n_2 and $f(n_2)$, have the same source and target cliques. \square

Weak summary size. The size of the weak summary is at most that of the original graph. This bound is reached if all nodes have distinct source and target properties. Moreover, Property 6 entails that the number of data edges in G/W is exactly the number of distinct data properties in G . Thus, its number of data nodes is at most twice this number. The number of type triples in G/W is bound by the smallest between the number of type edges in G .

4.2 Typed weak summary

In this section, we introduce a variant of the W -summary, based on the simple idea that RDF type information, when available, can be used to consider equivalent resources having the same *set of types* (recall that an RDF node may have 0, 1 or more types, not necessarily related). This summary represents all G nodes having the exact same types together, while for untyped resources, it reverts to the quotient by the weak equivalence relation \equiv_w . Such type-conditioned summarization can be seen as applying two consecutive quotients to G . We start by introducing two helper notions:

Definition 16. (TYPE-BASED SUMMARY) *The type-based summary of G , denoted G/T , is the summary of G through the \equiv_T equivalence relation.*

G/T groups together typed resources which have the same non-empty set of types, and copies each untyped G node, since none of them is equivalent to any other node. We assume available an *injective* function C which, called on a non-empty set of class URIs from G , returns a fresh URI (not in G). Figure 10 illustrates a type-based summary; observe the nodes produced by calls to C .

Further, we introduce:

Definition 17. (U-WEAK EQUIVALENCE AND SUMMARY)

Two nodes in a graph G are untyped-weak equivalent, denoted $n_1 \equiv_{UW} n_2$, iff n_1 and n_2 have no type in G and $n_1 \equiv_w n_2$.

The untyped-weak summary of G , denoted G/UW , is its summary through \equiv_{UW} .

Untyped weak equivalence is the weak equivalence \equiv_w restricted to untyped resources only. The untyped-weak summary G/UW treats the G triples whose subject and object lack types exactly like G/W , and leaves all typed resources untouched. We can now define:

Definition 18. (TYPED WEAK SUMMARY) *The typed weak summary G/TW of an RDF graph G is the untyped-weak summary of the type-based summary of G , namely $(G/T)_{UW}$.*

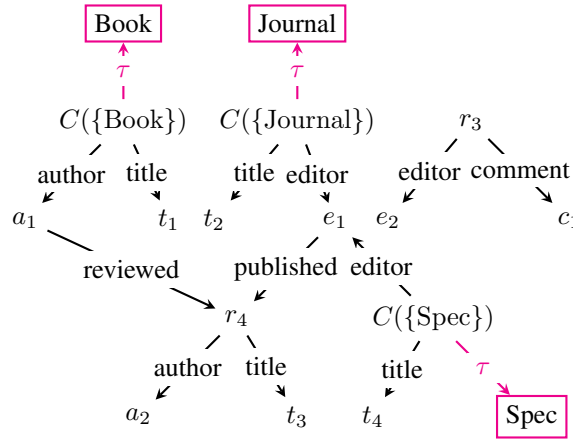


Figure 10: Type-based summary G_T of the RDF graph in Figure 3.

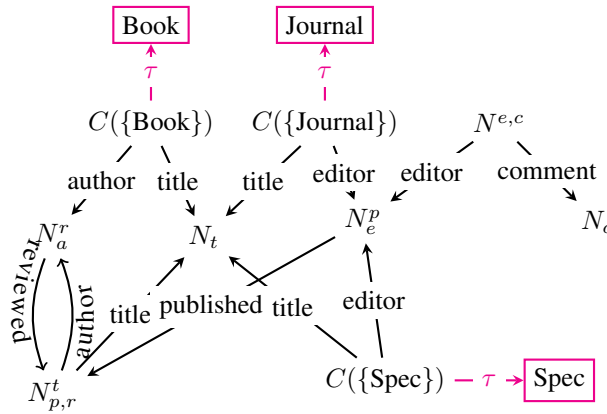


Figure 11: Typed weak summary G_{TW} of the graph G in Figure 3.

For example, Figure 11 shows the typed weak summary of the RDF graph in Figure 3. Compared with G_T from which it is built (Figure 10), the typed weak summary represents all editors with a single node, all titles with a single node etc., like G_W (Figure 8). However, in contrast with G_W , G_{TW} represents publications by three distinct nodes, depending on their types.

Proposition 8. (TW-INCOMPLETENESS) *Typed weak summaries are incomplete.*

Figure 12 shows a counter-example. In G and G_{TW} , resources are untyped; a typed resource only appears due to saturation and the constraint $a \leftrightarrow_d c$. Thus, in G_{TW} , one (untyped) node represents all data property subjects; in $(G_{TW})^\infty$, this (single) node gains a type, and in $((G_{TW})^\infty)_{TW}$, it is represented by a single node. In contrast, in G^∞ , one resource is typed and the other one isn't, leading to distinct nodes in $(G^\infty)_{TW}$. This cannot be isomorphic with the result obtained from G if one starts by TW summarization, thus by Definition 13, G_{TW} is not complete.

Typed weak summary size. The size of the typed weak summary is at most that of the original graph. This bound is reached if all nodes have distinct types.

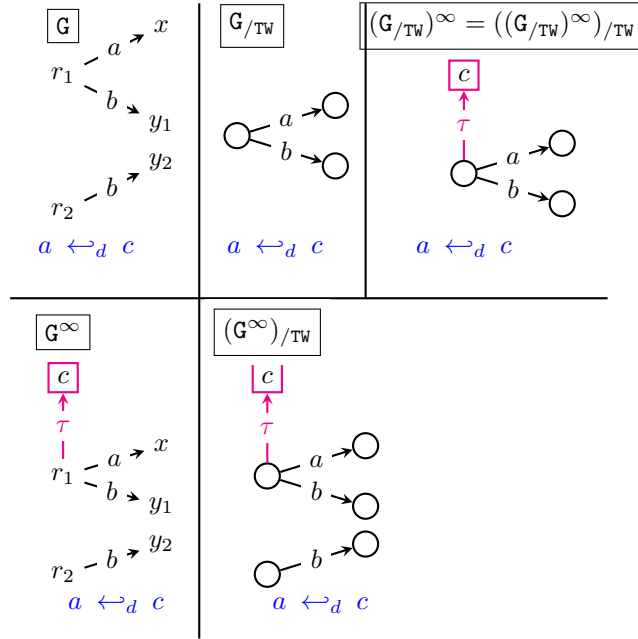


Figure 12: Typed weak summary completeness counter-example.

5 Strong equivalence summaries

We now present summaries based on the strong RDF node equivalence. Strong summaries (Section 5.1) are mainly based on this; typed strong summaries (Section 5.2) give preeminence to types.

5.1 Strong summary

Definition 19. (STRONG SUMMARY) *The strong summary S is an RDF summary based on the strong equivalence \equiv_s .*

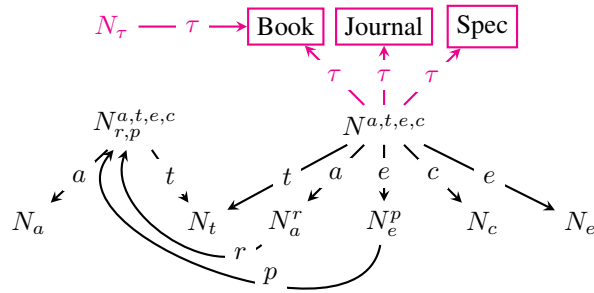


Figure 13: Strong summary of the RDF graph in Figure 3.

Recall from the definition of \equiv_s that only data nodes having the same source clique and the same target clique are equivalent in this sense. Thus, it is easy to establish a bijection between any pair (source

clique, target clique) of a data node in G , and a node in the strong summary. To follow through this intuition, in this section, we denote by N_{SC}^{TC} the node representing the set of \equiv_S -equivalent nodes of G having the target clique TC and the source clique SC .

For example, the strong summary of the graph of Figure 3 is shown in Figure 13. The strong summary comprises the nodes:

- $N(\emptyset, SC_1) = N^{a,t,e,c}$, for r_1, r_2, r_3 and r_5 ;
- $N(TC_5, SC_1) = N_{r,p}^{a,t,e,c}$, for r_4 ;
- $N(TC_1, SC_2) = N_r^a$, for a_1 ;
- $N(TC_2, \emptyset) = N_t$, for t_1, t_2, t_3 and t_4 ;
- $N(TC_3, SC_3) = N_e^p$, for e_1 ;
- $N(TC_3, \emptyset) = N_e$, for e_2 ;
- $N(TC_4, \emptyset) = N_c$, for c_1 ;
- $N(TC_1, \emptyset) = N_a$, for a_2 ;
- $N(\emptyset, \emptyset) = N_\tau$ for r_6 ;
- Book, Journal, and Spec are copied from the schema of G .

As a result, the number of data nodes in the data component D_{S_G} is bound by $\min(|D_G|_n, (|D_G|_e^0)^2)$ (recall the notations introduced in Section 2.1). Indeed, S_G cannot have more data nodes than G ; also, it cannot have more nodes than the number of source cliques times the number of target cliques, and each of these is upper-bounded by $|D_G|_e^0$. By a similar reasoning, the number of data triples in S_G is bound by $\min(|D_G|_e, (|D_G|_e^0)^4)$. In the worst case, T_{S_G} has as many nodes (and triples) as T_G ; S_{S_G} is identical to S_G .

Similarly to weak summaries, strong summaries are *complete* (Definition 13); this is our last major theoretical result. We prove this based on counterparts of statements established for G_W :

Proposition 9. (SAME CLIQUES-S) G^∞ and $(G/S)^\infty$ have identical source clique sets, and identical target cliques sets. Further, any node n in G^∞ has the source clique SC and target clique TC iff $f(n)$ in $(G/S)^\infty$ has the source clique SC and target clique TC , with f the representation function (recall Proposition 5).

The above Proposition follows directly from Proposition 5 and the following Lemma:

Lemma 3. Data properties are target-related (resp. source-related) in $(G/S)^\infty$ iff they are target-related (resp. source-related) in G^∞ .

Proof. “If”: if data properties are target-related in G^∞ , then they belong to a same target clique TC^∞ in G^∞ . Let n_1, \dots, n_k be the set of all G resources which are values of some properties in TC^∞ . By definition of an RDF summary and Proposition 5, each image of $f(n_i) = N(TC(n_i), SC(n_i))$ of n_i , for $1 \leq i \leq k$ is at least the object of the same properties as n_i , hence all the properties of TC^∞ in G^∞ are target-related in $(G/S)^\infty$.

“Only If”: if two data properties p_1 and p_2 are target-related in $(G/S)^\infty$, then they belong to a same target clique $TC^{S,\infty}$, in which they are at distance $n \geq 0$, i.e., they are target-related because of a set $\bigcup_{i=0}^n \{r_{i+1}\}$ of nodes (recall Definition 7) which all have the target clique $TC^{S,\infty}$. In G/S , each such r_{i+1} has a target clique $TC_i^S \subseteq TC^{S,\infty}$, moreover each r_{i+1} results from a set of G nodes $n_{i+1}^j, j \geq 1$, which by definition of a strong RDF summary, have all the source clique TC_i^S . Hence, every such n_{i+1}^j node has target clique $TC^{S,\infty}$ in G^∞ (since G and G/S have the same schema), in which p_1 and p_2 are target related. \square

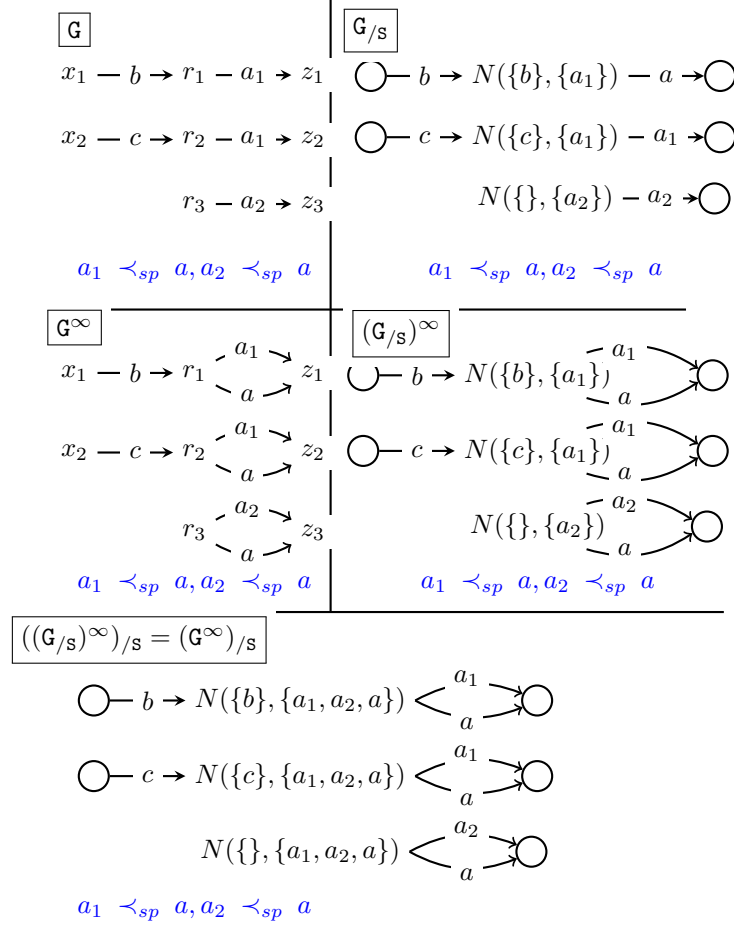


Figure 14: Illustration of the strong completeness statement (some summary nodes are shown as unlabeled circles).

Based on Proposition 9, we prove that G_S enjoys the sufficient condition stated by Theorem 1:

Theorem 4 (S-completeness). *Strong summaries are complete.*

Proof. We show that S summaries enjoy the sufficient condition for completeness stated in Theorem 1. That is, given two nodes n_1, n_2 in G^∞ , f the representation homomorphism corresponding to the weak-equivalence relation \equiv_S , and $f(n_1), f(n_2)$ the images of n_1, n_2 in $(G/S)^\infty$ through f (recall Proposition 5), let us prove that: $n_1 \equiv_S n_2$ in G^∞ iff $f(n_1) \equiv_S f(n_2)$ in $(G/S)^\infty$.

Two summary nodes are strongly-equivalent if they have the same non-empty source or target clique. The claim to prove hence immediately follows from Property 7, which guarantees that n_1 and $f(n_1)$, resp. n_2 and $f(n_2)$, have the same source and target cliques. \square

Figure 14 illustrates completeness on an example.

Strong summary size. The number of data nodes in the data component D_{S_e} is bound by $\min(|D_G|_n, (|D_G|_e^0)^2)$ (recall the notations introduced in Section 2.1). Indeed, S_G cannot have more data nodes than G ; also, it cannot have more nodes than the number of source cliques times the number of target cliques, and each

of these is upper-bounded by $|D_G|^0$. By a similar reasoning, the number of data triples in S_G is bound by $\min(|D_G|_e, (|D_G|^0)^4)$. In the worst case, T_{S_G} has as many nodes (and triples) as T_G ; S_{S_G} is identical to S_G .

5.2 Typed strong summary

The summary presented here is the counterpart of the typed weak one from Section 4.2, but based on strong equivalence.

Definition 20. (U-STRONG EQUIVALENCE AND SUMMARY) *Two nodes in a graph G are untyped-strong equivalent, denoted $n_1 \equiv_{US} n_2$, iff n_1 and n_2 have no type in G and $n_1 \equiv_S n_2$. The untyped-strong summary of G , denoted $G_{/US}$, is its summary through \equiv_{US} .*

Untyped strong equivalence restricts strong equivalence to untyped resources only. The summary $G_{/US}$ summarizes untyped data resources in strong fashion, and leaves all typed resources untouched (each is \equiv_{US} to itself). We can now define:

Definition 21. (TYPED STRONG SUMMARY) *The typed strong summary $G_{/TS}$ of an RDF graph G is the untyped-strong summary of the type-based summary of G , namely: $(G_{/T})_{/US}$.*

In our example, it turns out that $G_{/TS}$ for the RDF graph in Figure 3 coincides with the $G_{/TW}$ one shown in Figure 11. As can be easily seen from their definitions, the type-weak and type-strong summaries behave identically on the triples involving typed resources; on the untyped ones, the difference is of the same nature as the difference between the strong and weak summaries.

Proposition 10. (TYPED STRONG INCOMPLETENESS)

Typed strong summaries are incomplete.

More generally, in the presence of domain (\leftrightarrow_d) or range (\leftrightarrow_r) RDF schema constraints, one cannot compute $(G^\infty)_{/TS}$ from G because the TS summary represents typed resources differently from the untyped ones. The \leftrightarrow_d and \leftrightarrow_r constraints may turn untyped resources into typed ones, thus leading to divergent representations of the data nodes of G in $G_{/TS}$ and respectively $(G^\infty)_{/TS}$. Thus, one cannot build from G , without saturating it, a graph isomorphic to $(G^\infty)_{/TS}$.

Typed strong summary size. Similarly to the type weak summary, the size is at most that of the original graph. This bound is also reached if all nodes have distinct types.

6 Summary comparison

How do our summaries relate to each other, and how can one compare them? $G_{/W}$ and $G_{/S}$ are complete (Theorems 3 and 4), while $G_{/TW}$ and $G_{/TS}$ are not (Propositions 8 and 10).

Composing summaries. One can easily prove that $(G_{/S})_{/W} = G_{/W}$, i.e., one could compute $G_{/W}$ by first summarizing G into $G_{/S}$, and then applying weak summarization on this (typically much smaller) graph; similarly, one can show that $(G_{/TS})_{/TW} = G_{/TW}$. It is also the case that $(G_{/W})_{/S} = G_{/W}$, i.e., strong summarization cannot compress a weak summary further, and similarly $(G_{/TW})_{/TS} = G_{/TW}$. Figure 15 summarizes the main relationships between summaries.

Size comparison. As a consequence of the above, $G_{/W}$ has at most as many nodes and as most as many edges as $G_{/S}$ (weak summarizes at least as much as strong), and similarly, the typed weak summary $G_{/TW}$ summarizes at least as much as the typed strong summary $G_{/TS}$.

Accuracy comparison. Another interesting measure is the closeness between a graph and its summary. By Definition 11, any subgraph of the G can be embedded through homomorphism in any summary of

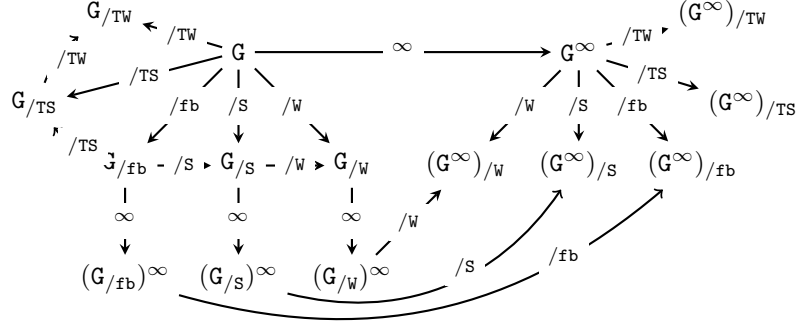


Figure 15: Main relations between our summaries.

G ; however, the opposite is not true, as the representation of data nodes by summary nodes may lead to subgraphs in the summary which do not exist in G . For instance, G/W in Figure 8 a node has outgoing edges labeled a and e , whereas such a node does not exist in the original graph (Figure 3). Formally:

Definition 22. (ACCURACY LOSS) *The accuracy loss of a summary $G_{/\equiv}$ w.r.t. G , denoted $AL(G_{/\equiv}, G)$, is the percentage of connected subgraphs of $G_{/\equiv}$ which are not isomorphic to any subgraph of G .*

The k -level accuracy loss, denoted $AL_k(G_{/\equiv}, G)$, for some $k \geq 2$, is the percentage of connected subgraphs of $G_{/\equiv}$ of size at most k , which lack isomorphic counterparts in G .

The lower AL , the more accurate the summary is w.r.t. G . Computing $AL(G_{/\equiv}, G)$ may be costly because of an excessively high number of summary subgraphs to check. Since $AL_k(G_{/\equiv}, G) \leq AL(G_{/\equiv}, G)$, we measure $AL_k(G_{/\equiv}, G)$ and use it as a lower bound approximation for AL . Note that small $G_{/\equiv}$ subgraphs show up in queries at least as frequently as larger ones, thus accuracy errors on small subgraphs is an interesting measure.

It is easy to see that for any graph G , summary $G_{/\equiv}$ and integers k_1, k_2 such that $k_1 < k_2$, $AL_{k_1}(G_{/\equiv}, G) \leq AL_{k_2}(G_{/\equiv}, G)$. This is because any summary subgraph of size k_1 which has no isomorphic image in the graph, is a subgraph of many subgraphs of size k_2 , which obviously are not isomorphic to a subgraph of G , either.

Clearly, AL and AL_k depend the RDF equivalence \equiv . At one extreme, if every node is only equivalent to itself, the summary has the size of G , and no accuracy loss. Conversely, if all data nodes are equivalent, the summary has a single node, having all types from G and an edge going from the node to itself for every data property in G . This maximizes AL for a given graph.

Relation with bisimulation summaries. Recast in our framework, the three RDF bisimulation summaries are complete (Theorem 2). Our W and S summaries are symmetric w.r.t. G edge source and target nodes, thus the comparison with the unidirectional FW and BW bisimulation summaries is limited. A deeper comparison is possible and interesting with the (symetric) FB -bisimulation summary $G_{/fb}$. Since $G_{/fb}$ preserves G structure and thus node cliques, we have $(G_{/fb})/W = G/W$ and $(G/W)/fb = G/W$, and similar statements for G/S , G/TW , and G/TS . Different from our summaries, $G_{/fb}$ guarantees $AL = 0$, but it is usually *much* larger, as illustrated in Section 3.2 and in our experiments, making it impractical for our goals.

7 Summarization algorithms

We have implemented a summarization tool (15.000 lines, Java 1.8) which, given an RDF graph G , builds G/W , G/TW , G/S and G/TS . G is initially stored in a PostgreSQL server, in a dictionary-encoded triple table; summaries are built mostly in memory and relying also on a few SQL queries. For efficiency, the algorithms for building G/W and G/TW are significantly different from following the summary definitions, and

we outline them below. G_S and G_{TS} are built in a fashion closer to the definitions. We introduce the data structures and then describe our algorithms.

| | Weak | Typed Weak | Strong | Typed Strong |
|----------------------------------|------|------------|--------|--------------|
| <u>rd</u> , <u>dr</u> | ✓ | ✓ | ✓ | ✓ |
| <u>dcls</u> | ✓ | ✓ | ✓ | ✓ |
| <u>clsd</u> | | ✓ | | ✓ |
| <u>dpSrc</u> , <u>dpTarg</u> | ✓ | ✓ | | |
| <u>srcDps</u> , <u>targDps</u> | ✓ | ✓ | | |
| <u>ntp</u> | ✓ | ✓ | ✓ | ✓ |
| <u>sc</u> , <u>tc</u> | | | ✓ | ✓ |
| <u>sToSc</u> , <u>oToTc</u> | | | ✓ | ✓ |
| <u>scToSrc</u> , <u>tcToTarg</u> | | | ✓ | ✓ |

Figure 16: Overview of the data structures for different summaries.

7.1 Data structures

Currently, our summary graphs are built in memory, based on the Trove library⁴ providing efficient collection data structures. We represent G URIs by integers, and store them in a set of map/ multi-map data structures, listed in Figure 16, together with the summaries which use them. All structures are not needed by all summaries, but many of them are used repeatedly. We describe them below, underlined for visibility.

rd maps G data nodes to their the G_{\equiv} representative; the dr multimap stores the reverse mapping (summary node to all G node it represents). We call *class set* a set of types attached to a same resource in G ; to each non-empty class set corresponds exactly one node in $G_{/TW}$ and $G_{/TS}$. clsd associates the corresponding summary node to each class set, while dcls stores the class set of every (typed) data node from G . In a weak summary G_w , every data property appears only once (Property 6); in $G_{/TW}$, appears *at most once with an untyped source and untyped target* (it may appear several times e.g., with subjects of different types and/or with typed vs. untyped subjects, see the four *title* edges in the $G_{/TW}$ summary in Figure 11). The maps dcls, dpSrc and dpTarg associate to a data property, its source and target (in G_w), respectively, its (possible) single untyped source and untyped target (in $G_{/TW}$). Conversely, for each G data node (for G_w), respectively, for each *untyped* G data node ($G_{/TW}$), the sets of data properties it is a source, respectively target of are stored in the maps srcDps (targDps). ntp maps every data property to the summary triple(s) it appears in. Each source and target clique is assigned an integer ID, and the IDs are stored in two lists, sc and tc. sToSc and oToTc associate to each data node, the ID of its source respectively target clique. Finally, for each source (target) clique ID, the G_S (and $G_{/TS}$) nodes having that clique are stored in scToSrc and tcToTarg.

7.2 Summarization algorithms

Building G_w and $G_{/TW}$ requires a single pass over the data. We create candidate summary nodes as we go over the G triples, and record the association between G nodes and the respective summary nodes in the maps described above, gradually gaining knowledge about the data properties that each G node is a source and target of. When nodes share a source or target property, their summary representative nodes are fused (again based on Property 6). Building G_S and $G_{/TS}$ require a first pass to compute the cliques, and a second one to build a summary node for each (source clique, target clique) pair associated to one or several G nodes.

⁴<http://trove.starlight-systems.com/>

Algorithm 1 Summarizing data triples in G/W **Input:** Data triples table D_G , the summary G/W **Output:** Data triples represented in G/W

```

1: for each  $s \ p \ o$  in  $D_G$  do
2:    $src \leftarrow \text{getSource}(s, p, G/W)$ 
3:    $targ \leftarrow \text{getTarget}(o, p, G/W)$ 
4:   //  $\text{getTarget}$  may have modified  $src$  and vice-versa
5:    $src \leftarrow \text{getSource}(s, p, G/W)$ 
6:    $targ \leftarrow \text{getTarget}(o, p, G/W)$ 
7:   if  $\text{existsDataTriple}(G/W, src, p, targ)$  then
8:      $\text{createDataTriple}(G/W, src, p, targ)$ 
9:   procedure  $\text{CREATEDATATRIPLE}(G/W, src, p, targ)$ 
10:     $\text{dtp.put}(p, src \ p \ targ)$ 
11:     $\text{dpSrc.put}(p, src), \text{srcDps.put}(src, p)$ 
12:     $\text{dpTarg.put}(p, targ), \text{targDps.put}(targ, p)$ 

```

Summarizing data triples in G/W . Algorithm 1 shows the procedure for summarizing data triples in a weak summary. The methods getSource and getTarget implement the representation functions, which map data nodes from G to data nodes in G/W . These methods create the respective nodes the first time they are called for a specific data property; on subsequent calls, the respective nodes are returned.

Algorithm 2 shows how we map the subjects of data triples in G to data nodes in G/W . After the method getSource has been executed, the source node of the property p , denoted src_u , and the node representing the subject s , denoted src_s must be the same. If neither src_u nor src_s exist yet, src will be a new data node representing s (line 5). In the cases when one of the nodes exist and the other does not, we simply use the existing node as src (lines 5-10). Moreover, if src_s had not existed, it means that s was unrepresented, so we assign src as its representative (line 7). On the other hand, if both src_u and src_s exist and they are the same, it does not matter which one we choose as src (line 12). When they differ, we have to merge them (line 14). The mergeDataNodes method replaces the node with less edges. The remaining node becomes the source/target of all the edges of the replaced node, and it is assigned to represent all the resources represented by the replaced node. Therefore, this method updates the replaced node in any of the maps $rd, dr, dpSrc, srcDps, dpTarg, targDps$, with the remaining node. Effectively, merging data nodes that are attached to common properties *gradually builds property cliques*. The method getTarget in Algorithm 1 is very similar to getSource , the only difference being in passing the object o instead of the subject, and working with untyped property targets instead of sources.

Algorithm 3 outlines the summarization of type triples in G/W . It iterates over the subjects and classes of all type triples and tries to represent each resulting pair as follows. We look up the summary data node src representing s . If s is attached to any data property in G , this means s has already been represented when summarizing data triples and we assign its representative data node to src and we add the class to the class set of src . Otherwise, s has some types but no data property, thus we record it as a typed-only resource. Once all type triples whose subjects also had some data property have been represented, if there are any typed-only resources, we represent them as well. The procedure $\text{representTypedOnly}$ creates the single data node N_τ representing all resources from the list of typed-only resources and having *all* the types of typed-only resources (recall Section 4.1).

7.2.1 Typed weak summary

Summarizing data triples in G/TW . The basic procedure for summarizing data triples in G/TW (Algorithm 4) is very similar to the one of weak presented in Algorithm 1, with the difference that the entries

Algorithm 2 Representing the subject s of a data property p in G with a data node in $G_{/w}$

Input: $s, p, G_{/w}$

Output: Data triples represented in $G_{/w}$ Variables:

src_u data node representing an untyped source of p
 src_s data node representing a (possibly typed) resource s

```

1: procedure GETSOURCE( $s, p, G_{/w}$ )
2:    $src_u \leftarrow dpSrc.get(p)$ 
3:    $src_s \leftarrow rd.get(s)$ 
4:   if  $src_u = \perp \wedge src_s = \perp$  then
5:     return createDataNode( $G_{/w}, s$ )
6:   else if  $src_u \neq \perp \wedge src_s = \perp$  then
7:      $rd.put(s, src), dr.put(src, s)$ 
8:     return  $src_u$ 
9:   else if  $src_u = \perp \wedge src_s \neq \perp$  then
10:    return  $src_s$ 
11:  else if  $src_s = src_u$  then
12:    return  $src_s$ 
13:  else
14:    return mergeDataNodes( $G_{/w}, src_s, src_u$ )
15: procedure CREATEDATANODE( $G_{/w}, r$ )
16:   $d \leftarrow newInteger$ 
17:   $rd.put(r, d), dr.put(d, r)$ 
18:  return  $d$ 

```

are stored to the maps only if src and $targ$ respectively, are untyped⁵. Further, the mapping of data nodes in $G_{/TW}$ is similar to the one of weak summary; they diverge when both src_s and src_u exist. *The decision on merging data nodes depends on the typing of src_s .* If src_s is typed, or it is untyped but already the same as src_u , we choose src_s as src (lines 12-13). On the other hand, if the two nodes are both untyped and different from each other, we must merge them (line 15). Therefore, in the typed weak summary only untyped data nodes may be merged, while the weak allows for the merging of typed data nodes as well. More precisely, since in the weak summary the data triples are represented *before* the type triples, the typing information is not yet available - the merged nodes may or may not be typed. However, in the typed weak summary the type triples are represented first and taken into account during the summarization of data triples to build a finer-grained summary.

Summarizing type triples in $G_{/TW}$. Algorithm 5 shows how we summarize type triples in the typed weak summary. For each distinct subject we retrieve its class set (line 3) and the representative data node of the class set (line 4). Observe that in a typed weak summary there is one data node *per distinct class set* which represents all subjects having the class set. If such a data node already exists for the given class set, we simply store it as the representative of the current subject (line 6). Otherwise, we create a new data node d , which represents s and has the same class set as s (line 8). Further, d is stored as the representative of its class set (line 9).

Summarizing schema. Prior to summarization, the encoded schema table was created. Since the schema of G and $G_{/w}$ is the same, no action is needed. **JDBC.** The results of all queries issued to Postgres are fetched through JDBC. A JDBC parameter that can have an impact on the summarization time is the fetch

⁵Since in $G_{/TW}$ all typed data nodes are created before untyped ones and the assigned values are sequential, the method *isTyped(d)* is as simple as checking whether d is less than or equal to the highest typed data node.

Algorithm 3 Summarizing type triples in G/W **Input:** Type triples table T_G , the summary G/W **Output:** Type triples represented in G/W

```

1:  $toRes \leftarrow []$ ;  $toCls \leftarrow []$ 
2: for each  $(s, c)$  in  $T_G$  do
3:    $repr \leftarrow \text{representTypeTriple}_{G/W}(s, c)$ 
4:   if  $repr = false$  then
5:      $toRes.add(s)$ ,  $toCls.add(c)$ 
6: if  $toRes.size > 0$  then
7:    $\text{representTypedOnly}(G/W, toRes, toCls)$ 
8: procedure REPRESENTTYPE TRIPLE( $G/W, s, c$ )
9:    $d \leftarrow rd.get(s)$ 
10:  if  $d = \perp$  then
11:    return false
12:   $cls_d \leftarrow dcls.get(d)$ ,  $cls_d.add(d, c)$ 
13:  return true
14: procedure REPRESENTTYPEDONLY( $G/W, toRes, toCls$ )
15:   $N_\tau \leftarrow \text{newInteger}$ 
16:  for each  $r$  in  $toRes$  do
17:     $rd.put(r, N_\tau)$ ,  $dr.put(N_\tau, r)$ 
18:   $cls_d \leftarrow dcls.get(N_\tau)$ 
19:  for each  $c$  in  $toCls$  do
20:     $cls_d.add(c)$ 

```

size. Larger fetch size is better because there are less real connections to the database to get all the results. But too large a fetch size will make you wait until that many results are available. The optimal fetch size will vary depending on the hardware setting.

7.2.2 Strong summary

Summarizing data triples in G/S . Instead of building the cliques gradually as in the weak summary, in Algorithm 6 we compute all source and target cliques at the very beginning (lines 1-2). Then in lines 3-15, for each data node d in G , a representative data node d_{repr} is assigned in G/S . This data node d_{repr} must have the exact same source and target clique as d , and there can be at most one such node for any two source and target cliques. Finally, for each data triple $s p o$ in D_G , we retrieve the representative data node of s , denoted src , and of o , denoted $targ$, forming a data triple $src p targ$ in G/S . The method for creating the data triple is the same as for the weak summary in Algorithm 1.

The procedure $\text{buildSourceCliques}(D_G, G/S)$ invokes the procedure $\text{computeSourceClique}(G/S, s, P_s)$ shown in Algorithm 7 for each subject in D_G and the set of data properties of which it is the source.

First, in Algorithm 7 (lines 3-6) we try to find an existing source clique, denoted $srcClq$, that shares at least one property with P_s . If such a clique exists we simply add to it all the properties from P_s (line 8). Otherwise, $srcClq$ is a new set with all the properties of P_s , and we add it to the list of all source cliques sc in G/S (lines 10-11). Finally, $srcClq$ is assigned to s in $sToSc$ (line 12).

The procedures for building target cliques are very similar to these described for the source cliques, so we omit the discussion.

Summarizing type triples in G/S . The types in G/S are summarized in the same manner as for the W_G , described in Algorithm 3.

Algorithm 4 Representing the subject of a data property in G with a data node in $G_{/TW}$ **Input:** The subject s of the data triple, the property p of s , the summary $G_{/TW}$ **Output:** Data triples represented in $G_{/TW}$

```

1: procedure GETSOURCE( $s, p, G_{/TW}$ )
2:    $src_u \leftarrow dpSrc.get(p)$ 
3:    $src_s \leftarrow rd.get(s)$ 
4:   if  $src_u = \perp \wedge src_s = \perp$  then
5:     return createDataNode( $G_{/TW}, s$ )
6:   else if  $src_u \neq \perp \wedge src_s = \perp$  then
7:      $rd.put(s, src), dr.put(src, s)$ 
8:     return  $src_u$ 
9:   else if  $src_u = \perp \wedge src_s \neq \perp$  then
10:    return  $src_s$ 
11:  else if  $src_u \neq \perp \wedge src_s \neq \perp$  then
12:    if isTyped( $src_s$ )  $\vee (src_s = src_u)$  then
13:      return  $src_s$ 
14:    else
15:      return mergeDataNodes( $G_{/TW}, src_s, src_u$ )

```

Algorithm 5 Summarizing type triples in $G_{/TW}$ **Input:** Type triples table T_G , the summary $G_{/TW}$ **Output:** Type triples represented in $G_{/TW}$

```

1:  $typ \leftarrow evalSELECT\ s, c\ FROM\ T_G\ ORDER\ BY\ s$ 
2: for each distinct  $s$  in  $typ$  do
3:    $cls_s \leftarrow dcls.get(s)$ 
4:    $d \leftarrow clsd.get(cls_s)$ 
5:   if  $d \neq \perp$  then
6:      $rd.put(s, d), dr.put(d, s)$ 
7:   else
8:      $d \leftarrow createDataNodeG_{/TW}, cls_s, s$ 
9:      $clsd.put(cls_s, d)$ 

10: procedure CREATEDATANODE( $G_{/TW}, r, cls_r$ )
11:   $d \leftarrow newInteger$ 
12:   $rd.put(r, d), dr.put(d, r), dcls.put(d, cls_r)$ 
13:  return  $d$ 

```

Algorithm 6 Summarizing data triples in G/S **Input:** Data triples table D_G , the summary G/S **Output:** Data triples represented in G/S

Variables:

 $scId, tcId$ - source (target) clique ID D_{SC}, D_{TC} - the set of source (target) data nodes in G/S representing the source clique $scId$ (target clique $tcId$)

```

1: buildSourceCliques $D_G, G/S$ 
2: buildTargetCliques $D_G, G/S$ 
3: for each distinct data node  $d$  in  $D_G$  do
4:    $scId \leftarrow sToSc.get(d), tcId \leftarrow oToTc.get(d)$ 
5:    $D_{SC} \leftarrow scToSrc.get(scId)$ 
6:    $D_{TC} \leftarrow tcToTarg.get(tcId)$ 
7:   if  $D_{SC} \cap D_{TC} \neq \emptyset$  then
8:      $rd.put(r, d_{repr}), dr.put(d_{repr}, r)$ 
9:   else
10:     $d_{repr} \leftarrow createDataNode(G/S, d)$ 
11:     $scToSrc.put(scId, d_{repr})$ 
12:     $tcToTarg.put(tcId, d_{repr})$ 
13: for each  $s\ p\ o$  in  $D_G$  do
14:    $src \leftarrow rd.get(s), targ \leftarrow rd.get(o)$ 
15:   if  $existsDataTriple(G/S, src, p, targ)$  then
16:      $createDataTriple(G/S, src, p, targ)$ 

```

Algorithm 7 Computing source cliques in G/S **Input:** The summary G/S , subject s from D_G , the set of data properties P_s of which s is the source**Output:** sc list and $sToSc$ map of G/S populated

```

1: procedure COMPUTESOURCECLIQUE( $G/S, s, P_s$ )
2:    $srcClq \leftarrow \perp$ 
3:   for each  $c$  in  $sc$  do
4:     if  $c \cap P_s \neq \emptyset$  then
5:        $srcClq \leftarrow c$ 
6:       break;
7:   if  $srcClq \neq \perp$  then
8:      $srcClq.addAll(P_s)$ 
9:   else
10:     $srcClq.addAll(P_s)$ 
11:     $sc.add(srcClq)$ 
12:    $sToSc.put(s, sc.indexOfsrcClq+1)$ 

```

| Graph G | $ G $ | $ G/w $ cf_w | $ G/Tw $ cf_{Tw} | $ G/s $ cf_s | $ G/TS $ cf_{TS} |
|--------------------------|-------|----------------------|-----------------------|----------------------|-----------------------|
| John Peel (BBC Radio) | 271k | 33 8223 | 51 5320 | 99 2741 | 61 4448 |
| INSEE Geo | 369k | 311 1186 | 360 1025 | 317 1164 | 361 1022 |
| DBpedia Person | 8M | 10 788957 | 12 657464 | 28 281770 | 17 464092 |
| DBLP | 150M | 71 2123767 | 258 584447 | 206 731978 | 285 529079 |
| LinkedCT | 49M | 442 111050 | 505 97196 | 729 67331 | 598 82081 |
| LinkedMDB | 6.14M | 8 118155 | 12 78770 | 8 118155 | 12 78770 |
| LUBM | 1M | 162 7579 | 267 4599 | 208 5903 | 267 4599 |
| LUBM | 10M | 162 74013 | 267 44907 | 206 58204 | 267 44907 |
| LUBM | 100M | 162 705897 | 267 428297 | 209 547155 | 267 428297 |
| WatDiv | 10M | 124 88035 | 982 11117 | 249 43841 | 1187 9196 |
| BSBM | 1M | 1046 956 | 3054 327 | 1067 937 | 3054 327 |
| BSBM | 10M | 4066 2591 | 14813 711 | 4101 2569 | 14813 711 |
| BSBM | 100M | 13390 7775 | 40270 2585 | 13429 7753 | 40270 2585 |

Figure 17: Summary graph sizes and compression factors (rounded down to the closest smaller integer) for various graphs.

Proposition 11. (ALGORITHM CORRECTNESS) *Given an RDF graph G: Algorithms 1 and 3, together, build the weak summary G/w . Algorithms 4 and 5 build the typed weak summary G/Tw . Algorithm 6 shows how to build the strong summary S_G ; building TS_G is very similar.*

Algorithm complexity. The complexity of building G/w and G/Tw is dominated by the cost of merging summary nodes, in the worst case, once per triple in G, leading to $O(|G| \cdot DPG)$, where DPG is the number of distinct data properties in G. To build G/s and G/TS , we actually compute the cliques. First, we find the distinct data properties of every distinct data node in G, at a cost of $O(N \cdot \log(N) + |G|)$, where N is the number of data nodes in G. Then, for each data node n , we examine all the cliques known so far, to see if some should be fused due to the data properties of n : in the worst case we may do DPG fusions at a cost of DPG each, thus a total cost of DPG^2 . Further, if there are fusions, the associations between the data nodes previously visited and their cliques need to be updated. A (very pessimistic) upper bound for the number of nodes concerned is N. Thus, the complexity of this stage is in $O(N(N + DPG^2))$; this dominates the complexity of the first step and thus dictates the complexity of building G/s and G/TS . As our experiments show next, the actual scale-up is much better, and close to linear in $|G|$.

Computing accuracy loss. We developed a tool which enumerates summary subgraphs and for each subgraph, asks an existential SQL query to check if the subgraph has matches in the data. Our algorithm enumerates smallest subgraphs first, and re-uses the knowledge that a certain subgraph has no matches in G to avoid asking queries that correspond to larger graphs comprising the empty one (as we can ascertain the larger graph has no matches, either).

8 Experimental results

Settings. We implemented our algorithms in Java 1.8, and relied on Postgres for storing the integer-encoded triple table. The computer had an Intel Xeon CPU E5-2640 v4 @2.40GHz and 124 GB of memory; we used PostgreSQL v9.6 with 30 GB of shared buffers and 640 MB working memory. The JVM had 70 GB of RAM.

Datasets. We summarized *synthetic* RDF graphs: BSBM [3], LUBM [12] and WatDiv (<http://dsg.uwaterloo.ca/watdiv/>) benchmark datasets, as well as *real-life* graphs: the John Peel BBC Radio Sessions (<http://dbtune.org/bbc/peel/>), French INSEE Geo geographic data (<http://rdf.insee.fr/geo/index.html>), DBLP, the DBpedia Person dataset (<http://wiki.dbpedia.org/Downloads2015-04>), Linked Clinical Trials (LCT) [24] and LinkedMDB [14]. Their sizes (in triples) are denoted $|G|$ in Figure 17.

Summary size. We define the **compression factor** cf_{\equiv} as $|G|/|G_{\equiv}|$, the ratio between the numbers of triples in G and in the summary. (Since schema triples are the same, compression comes from the representation of G data triples through G_{\equiv} triples.) Figure 17 shows the size of various summary and the resulting compression factor underneath; cf values higher than 10^4 appear in bold font. The compression factor is at least 300, more usually above 10^3 , and up to $2 \cdot 10^6$ for DBLP dataset; G_w is typically smallest, while the G_{TW} and G_{TS} may be up to one order of magnitude larger. The LUBM and BSBM 1M and 100M triples synthetic datasets offer an interesting contrast. LUBM compresses rather well; the number of data edges does not increase as the data grows. This enables higher cf values from $4 \cdot 10^3$ to $7 \cdot 10^5$. BSBM is harder to compress, which correlates with its higher number of distinct data properties (642 to 8082), which also grows with the data size, classes grow from 159 to 2019, and schema triples from 150 to 2010). BSBM compression factors range from 327 to 7775, more modest but still very significant. *The compression factor increases with the data size* in both cases, highlighting summarization’s particular interest for large graphs. For LinkedMDB, compression factors higher than 10^5 contrast with **the compression factor of 1.42** reported in [19] for the **fb-bisimulation summary**.

A more detailed look at BSBM summarization is provided in Figure 18, which shows node and edge counts for summaries of BSBM graphs of various sizes. The horizontal axis is labeled in input triples; the vertical axis is in log scale. For reference, we plot the number of class nodes of G next to the summary data nodes, and the number of schema edges next to the edge counts. The figure shows that G_w and G_s strongly summarize graph structure, while G_{TW} and G_{TS} which isolate typed from untyped data nodes lead to larger, more complex summaries.

Summarization time. Figure 19 shows the average time (in seconds, 3 executions) to build summaries; they go up to 800 seconds for 100M triples. Times for the LUBM and BSBM synthetic data are plotted (both axis are logarithmic, the horizontal axis labeled in $|G|$ triples). The graphs show that summarization time grows almost linearly, and the time to build G_w grows fastest. This is because building G_w and G_s maximizes the number of nodes for which cliques are analyzed and fused; in contrast, when building G_{TW} and G_{TS} , *the cliques of typed data nodes are never constructed (directly or indirectly)*. We also see that G_s and G_{TS} are oftentimes more expensive to build than G_w and G_{TW} ; this is because of the extra effort needed to build the cliques, while G_w and G_{TW} use a more direct, faster method. We find summarization times acceptable, also considering that this is an off-line task.

Building complete summaries. Figure 20 shows the time to build the *complete* weak and strong summaries: by saturating and then summarizing ($(G^\infty)_w$ and $(G^\infty)_s$); and separately, by summarizing G first, then saturating and then summarizing again, taking advantage of Theorems 3 and 4. For saturation, we used the algorithm described in [10]. \times_w , respectively \times_s denote the speed-up factors achieved thanks to the completeness theorems. Figure 20 shows that computing $(G^\infty)_w$, $(G^\infty)_s$ by the indirect method given by the completeness theorems is much faster; indeed, this method reduces the size of the graph in

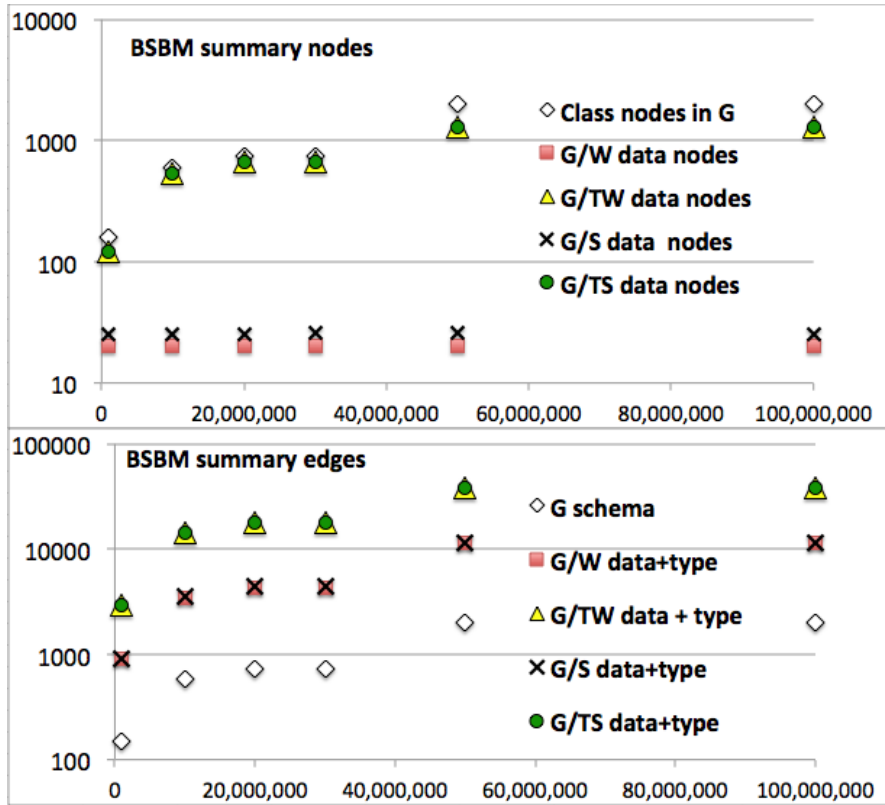


Figure 18: Summary node counts (top) and edge counts (bottom) for BSBM graphs. the first place, making subsequent operations more efficient.

| Dataset | Weak | Strong | Typed Weak | Typed Strong |
|-----------|----------------------------|----------------------------|----------------------------|----------------------------|
| John Peel | $AL_2 = 81$ $AL_3 = 96$ | $AL_2 = 29$ $AL_3 = 53$ | $AL_2 = 15$ $AL_3 = 30$ | $AL_2 = 14$ $AL_3 = 30$ |
| INSEE | $AL_2 = 41$ $AL_3 = 68$ | $AL_2 = 36$ $AL_3 = 64$ | $AL_2 = 1$ $AL_3 = 5$ | $AL_2 = 1$ $AL_3 = 5$ |
| LUBM 1M | $AL_2 = 61$ $AL_3 = 83$ | $AL_2 = 30$ $AL_3 = 53$ | $AL_2 = 0$ $AL_3 = 0$ | $AL_2 = 0$ $AL_3 = 0$ |
| DBpedia | $AL_2 = 0$ $AL_3 = 2$ | $AL_2 = 0$ $AL_3 = 4$ | $AL_2 = 0$ $AL_3 = 1$ | $AL_2 = 0$ $AL_3 = 0$ |
| DBLP | $AL_2 = 48$ | $AL_2 = 38$ | $AL_2 = 44$ | $AL_2 = 39$ |
| BSBM 1M | $AL_2 = 78$ | $AL_2 = 0$ | $AL_2 = 65$ | $AL_2 = 0$ |

Table 2: Bounded accuracy loss values.

Accuracy loss. Table 2 shows the accuracy loss AL_k (in %, rounded to the lower closest integer) for $k \in \{2, 3\}$ on various RDF graphs; the lowest value(s) on each row are in boldface. G/W and G/S have comparable accuracy losses, while on the other hand G/TW and G/TS are significantly more accurate. The difference is dramatic for LUBM and DBpedia, where TW and TS have 0 accuracy loss for $k = 2$, just as the DBpedia TS summary for $k = 3$, while the DBpedia TW summary has a minimum loss of 1% for $k = 3$.

Experiment conclusion. Our experiments have shown that our summaries can be built efficiently, and they reduce graph sizes by factors hundreds up to $2 \cdot 10^6$. This confirms their power of compactly

| Dataset | G | G/w | G/s | G/TW | G/TS |
|-----------|-------|--------|--------|--------|----------|
| John Peel | 271k | 0.33 | 1.37 | 0.43 | 0.98 |
| INSEE Geo | 369k | 0.58 | 1.65 | 0.55 | 1.50 |
| DBpedia | 8M | 9.73 | 25.85 | 6.82 | 21.80 |
| DBLP | 150M | 434.88 | 636.56 | 175.87 | 12412.49 |
| LinkedCT | 42M | 74.55 | 314.75 | 88.73 | 252.38 |
| LinkedMDB | 6.14M | 1.12 | 3.49 | 0.90 | 3.06 |
| LUBM | 1M | 1.34 | 3.30 | 1.33 | 3.20 |
| LUBM | 10M | 13.92 | 28.27 | 10.42 | 31.14 |
| LUBM | 100M | 221.62 | 403.97 | 132.89 | 457.63 |
| WatDiv | 10M | 14.78 | 22.77 | 8.33 | 21.02 |

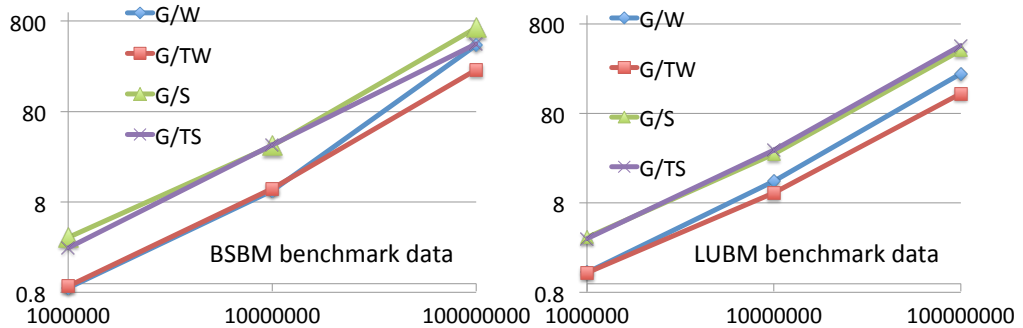


Figure 19: Summarization times (s) for various RDF datasets

| Dataset | $(G^\infty)/w$ | $((G/w)^\infty)/w$ | \times_w | $(G^\infty)/s$ | $((G/s)^\infty)/s$ | \times_s |
|-----------|----------------|--------------------|------------|----------------|--------------------|------------|
| INSEE Geo | 49.17 | 1.25 | 39.33 | 45.94 | 1.70 | 27.02 |
| LUBM 1M | 20.77 | 1.42 | 14.62 | 21.27 | 2.72 | 7.8 |
| LUBM 10M | 220.36 | 12.31 | 17.90 | 219.64 | 34.04 | 6.45 |
| BSBM 1M | 22.88 | 1.42 | 16.11 | 16.33 | 3.25 | 5.02 |
| BSBM 10M | 301.98 | 134.06 | 2.25 | 213.73 | 71.13 | 3.01 |
| DBLP | 465.09 | 14.83 | 31.36 | 434.87 | 43.28 | 10.04 |

Figure 20: Times (s) to generate semantically complete summaries.

representing RDF graphs. G/s and G/TS take longer to compute than G/w and G/TW since the former actually compute cliques, while the latter rely instead on Property 6 for a more efficient method. Weak summaries are often less accurate, but the opposite also occurs, e.g., on the DBpedia person dataset; weak summaries compress more (and thus are smaller). If summary size or building time are a concern, we recommend using G/w ; otherwise, G/s may give finer-granularity information.

9 Related Work

Graph summarization has a long history and many applications. We survey general graph summarization works (with no specific focus on RDF graphs) in Section 9.1, then discuss RDF-specific summarization proposals in Section 9.2.

9.1 Non-RDF graph summaries

Such summaries ignore graph semantics, thus summary completeness, a main focus of our work, is not in general guaranteed.

Several works built graph summaries as a support for graph indexing. In this context, for each summary node, we store the IDs of the original graph nodes represented by the summary node; this set is typically called the extent. Given a query, evaluating the query on the summary and then using the summary node extents allows to obtain the query results.

Dataguides [11], and in particular the so-called strong dataguide, were among the first proposals for structural summaries (also called a-posteriori schemas) of RDF graphs. Dataguides do not compare directly with our approach, since they are not quotient graphs: in particular, a graph node may be represented by more than one dataguide node, which is impossible in a quotient-based summary. A dataguide may be larger than the original graph, and its construction has worst-case exponential time complexity.

[25, 6, 16, 8] build quotient summaries based on (sometimes unidirectional, sometimes k -bounded) bisimulation; we have discussed the relationship between our and bisimulation summaries in Section 6. The complexity of building a bisimulation summary is typically $O(E \cdot \log(N))$, where E is the number of edges and N is the number of nodes in G .

[24] explores efficient answering of neighbor queries over a compressed social network graph, while minimizing the number of bits per edge required for storage, by relying on covers of an input graph. Thus, the focus is entirely different than ours, which is efficiently representing all RBGP queries from an RDF graph G into its summaries, while preserving the semantics of G , by grouping similar nodes based on equivalence classes.

[9] considers reachability and graph pattern queries on labeled graphs, and builds *answer-preserving summaries* for such queries; this property subsumes accuracy and representativeness. However, their query semantics is not based on graph homomorphisms, which underlie SPARQL and the dialects we use, but on *bounded graph simulation*. Under this semantics, answering a query becomes P (instead of NP), at the price of not preserving the query structure (i.e., joins).

[32] proposes the SNAP (Summarization on Grouping Nodes on Attributes and Pairwise Relationships) technique which, given a set of properties the user is interested in, produces a summary that groups the interesting nodes in the incoming RDF graph, based on their roles with respect to the user-specified properties. (The authors consider graphs where nodes have “attributes” and “relationships”; this can be seen as a restriction of RDF where we do not consider types nor schema information.) SNAP partitions nodes into classes that may end up being too many and too small (two few nodes will be equivalent). For more flexibility, the authors also propose k -SNAP where the similarity requirement between nodes of the same class is relaxed and the summary is computed taking into account a space budget (i.e., a maximum summary size). Thus, while SNAP is based on strict graph relations, k -SNAP mixes in a quantitative component (with a threshold of 50% proposed by the authors at some point etc.). Our work is different in that we request no user input and propose complete, logical-oriented summaries.

[5] introduces an aggregation framework for OLAP operations on labeled graphs. The authors assume available an OLAP-style set of dimension with their hierarchies, and measures, and based on them they define a “graph cube” and investigate efficient methods for computing it. The authors also frame graph topological information as dimensions, and use them for graph aggregation. This is quite different from our framework, which is based on logical representation of structure and semantics.

[4] considers node- and edge-labeled graphs and focuses on building a set of randomized summaries, so that one can mine the summary set instead of the original graph for better performance, with guaranteed bounds on the information loss thus incurred. Semantics and properties such as those of our summaries are out of their scope.

[39] seeks to identify frequent subgraphs and store extent information from those, so as to answer queries asking for the respective part of the graph. This approach by design does not reflect all data and

is based on numeric information, as opposed to our summaries. [40] is very similar, but considers trees instead of graphs.

A graph *core* C for a given graph G is a graph such that an isomorphism exists between G and C , and C is the smallest graph with this property⁶. Our summaries are not cores of G , since we cannot guarantee a homomorphism from any of them to the G . In exchange, all of our summaries can be built in polynomial time in the size of G , while computing the core is much harder.

A separate brand of works is interested in computing approximate summaries of graphs, to be used to efficiently answer queries with an acceptable level of error [26, 18]. The focus is on graph compression while preserving bounded-error query answering and/or the ability to fully reconstruct the graph from the summary, with the help of so-called “corrections” (edges to add to or remove from the “expansion” of the summary into the regular part of the graph it derives from). [22] focuses on summarizing undirected, node-labeled graphs; nodes are partitioned, while edge information is preserved under the form of the number of edges within each partition set, and number of edges between every two such sets. These works contrast with the full representativity, accuracy and completeness we aim at.

[23] surveys many other quantitative, mining-oriented graph sampling and summarization methods.

9.2 RDF graph summaries

Several works have more recently target the construction of RDF graph summaries. Still, most ignore RDF semantics, or make limiting assumptions on the input RDF graph, while we target general graphs, with or without schema statements. Overall, our work is the first to consider the issue of building complete RDF summaries, and to provide a method for doing it and a sufficient completeness condition.

[36] builds a tree-structured index of an RDF graph by picking some “centroid” RDF graph nodes and organizing information about the graph in a tree; if one views the index as a summary, it is not a quotient, and since inference is not considered, it is not complete, either.

[28] introduces a *simulation* RDF quotient based on *triple (not node) equivalence*; the summary is not an RDF graph. (Intuitively, *simulation* is a unidirectional “half” of bisimulation, e.g., to define FW RDF similarity, remove item (ii) from Definition 9). For compactness, they bound the simulation, for small k values, e.g., 2; this enables compression factors of about 10^4 . They use summaries to reduce query join effort, but do not consider the possible schema, nor summary completeness.

[30] considers the problem of efficiently building a quotient summary of RDF graphs based on the bisimulation relationship between nodes. The authors present several algorithms, in particular based on MapReduce, but do not reserve any special treatment to RDF class and/or property nodes; this would allow in particular to summarize different classes into a single node, e.g., if the classes only participate to triples as subjects of the subclassOf property. Such schema triple summarization compromises completeness, as it is not possible to build the summary of G^∞ out of a summary that does not preserve the schema of G .

[33] has as main goal to answer queries over RDF graphs and proposes to help this based on a structure index, which is defined as a bisimulation quotient; the authors show that the summary is representative. Further, the authors study limited versions of the bisimulation quotient by considering: only forward bisimulation; only backward bisimulation; only neighborhoods of a certain length. The authors do not consider graph semantics, nor completeness.

[29] presents SAP HANA’s approach for working with graph data. They consider a restricted RDF model where every node has exactly one type; to simplify the experience of users querying graphs, they propose that users specify some simple summarization rules (summarization here viewed as aggregation) together with their query, rules which the system applies to group the query results. Beyond its restricted data model and lack of support for RDFS inference, the work has clearly a very different focus and does not aim to produce an interesting representation of the whole graph as we do.

⁶[https://en.wikipedia.org/wiki/Core_\(graph_theory\)](https://en.wikipedia.org/wiki/Core_(graph_theory))

[21] introduces an RDF summary graph which they use to support keyword search in RDF graphs. They restrict RDF graphs so that (i) one resource can have at most one type, and (ii) same-type resources have similar sets of data properties. These restrictions make the summarization problem significantly different from ours, as we do not make such restrictive assumptions. The authors do not consider inference either.

[13] applies a graph clustering algorithm (in particular METIS⁷) to determine the best partitioning to use in order to store an RDF graph in a distributed system. The partitioning criterium is quantitative (METIS seeks to minimize edges across partitions, thus strong connections will tend to keep nodes together). Such partitioning is not a quotient graph and properties such as accuracy and completeness cannot be guaranteed.

[7] also considers summaries defined as quotients based on bisimulation, disregarding the role played by schema triples, and proposes a highly parallel implementation based on the GraphChi framework. They show how to use the summary as a support for query evaluation: incoming SPARQL queries are evaluated on the summary and then the results on the summary are transformed in results on the original graph by exploring the extents of summary nodes.

[17] uses a density-based clustering algorithm to identify nodes with similar structure in an RDF graph; the goal is to propose an RDF type for each cluster of resources and identify relations between these types. Thus, the authors do not consider RDF schema information that the graph may already have, nor are they interested in accuracy, completeness etc. for the result of their clustering.

[35] aims to summarize an RDF Schema based also on the frequency of representation of various schema features in the data graph. The approach makes several limitative assumptions on the RDF data model, e.g., that each resource has exactly one type etc. While our summaries do not reduce the size of an RDF Schema, they allow to significantly compress general RDF graphs (with or without schema information) and enjoy several interesting properties.

[31] focus on summarizing RDF graphs which may have subtype and subproperty triples, but they do not consider domain and range, and instead rely on *abstract knowledge patterns*, expressing something related but different: the pattern $c_1 \text{ p } c_2$ states that a dataset holds at least one resource of type c_1 having the property p whose value is a resource of type c_2 . (In contrast, recall that the domain constraint $p \leftrightarrow_d c_1$ states that any resource having property p is also of class c_1 .) A summary is a set of types, a set of AKPs, and some statistics; the authors are only concerned with the minimality of their summaries and do not consider properties such as accuracy, representativeness, completeness etc.

10 Conclusion

We studied the construction of *complete* summaries of RDF data, accounting for all its structure and semantics; in particular, we focused on *quotient* graph summaries, widely studied in the literature. We introduced four novel summaries, two of which (as well as bisimulation-based ones) are complete. All our summaries are representative and accurate w.r.t. a large and useful dialect of SPARQL; they can be built efficiently and achieve RDF graph compression factors of 10^2 to $2 \cdot 10^6$; the highest values are reached for real-life datasets. Future work will improve scalability by leveraging a massively parallel platform such as Spark. We are also interested in devising advanced visualizations of our RDF summaries.

References

- [1] RDF Summary project website: <https://team.inria.fr/cedar/projects/rdfsummary>, 2017.
- [2] S. Abiteboul, R. Hull, and V. Vianu. *Foundations of Databases*. Addison-Wesley, 1995.
- [3] C. Bizer and A. Schultz. The Berlin SPARQL Benchmark. *Int. J. Semantic Web Inf. Syst.*, 5(2):1–24, 2009.

⁷<http://glaros.dtc.umn.edu/gkhome/metis/metis/overview>

- [4] C. Chen, C. X. Lin, M. Fredrikson, M. Christodorescu, X. Yan, and J. Han. Mining graph patterns efficiently via randomized summaries. *PVLDB*, 2(1), 2009.
- [5] C. Chen, X. Yan, F. Zhu, J. Han, and P. S. Yu. Graph OLAP: towards online analytical processing on graphs. In *ICDM*, 2008.
- [6] Q. Chen, A. Lim, and K. W. Ong. $D(K)$ -index: An adaptive structural summary for graph-structured data. In *SIGMOD*, 2003.
- [7] M. P. Consens, V. Fionda, S. Khatchadourian, and G. Pirrò. S+EPPs: Construct and explore bisimulation summaries + optimize navigational queries; all on existing SPARQL systems (demonstration). *PVLDB*, 8(12), 2015.
- [8] M. P. Consens, R. J. Miller, F. Rizzolo, and A. A. Vaisman. Exploring XML web collections with DescribeX. *ACM TWeb*, 4(3), 2010.
- [9] W. Fan, J. Li, X. Wang, and Y. Wu. Query preserving graph compression. In *SIGMOD*, 2012.
- [10] F. Goasdoué, I. Manolescu, and A. Roatiş. Efficient query answering against dynamic RDF databases. In *EDBT*, 2013.
- [11] R. Goldman and J. Widom. Dataguides: Enabling query formulation and optimization in semistructured databases. In *VLDB*, 1997.
- [12] Y. Guo, Z. Pan, and J. Hefflin. LUBM: A benchmark for OWL knowledge base systems. *J. Web Sem.*, 3(2-3), 2005.
- [13] S. Gurajada, S. Seufert, I. Miliaraki, and M. Theobald. Using graph summarization for join-ahead pruning in a distributed RDF engine. In *SWIM*, 2014.
- [14] O. Hassanzadeh and M. P. Consens. Linked movie data base. In *LDOW*, 2009.
- [15] M. R. Henzinger, T. A. Henzinger, and P. W. Kopke. Computing simulations on finite and infinite graphs. In *FOCS*, 1995.
- [16] R. Kaushik, P. Bohannon, J. F. Naughton, and H. F. Korth. Covering indexes for branching path queries. In *SIGMOD*, 2002.
- [17] K. Kellou-Menouer and Z. Kedad. Schema discovery in RDF data sources. In *ER*, 2015.
- [18] K. Khan, W. Nawaz, and Y. Lee. Set-based approximate approach for lossless graph summarization. *Computing*, 97(12):1185–1207, 2015.
- [19] S. Khatchadourian and M. P. Consens. Constructing bisimulation summaries on a multi-core graph processing framework. In *GRADES*, 2015.
- [20] D. Lanti, M. Rezk, G. Xiao, and D. Calvanese. The NPD benchmark: Reality check for OBDA systems. In *EDBT*, 2015.
- [21] W. Le, F. Li, A. Kementsietsidis, and S. Duan. Scalable keyword search on large RDF data. *IEEE TKDE*, 26(11), 2014.
- [22] K. LeFevre and E. Terzi. GraSS: Graph structure summarization. In *SDM*, 2010.
- [23] S.-D. Lin, M.-Y. Yeh, and C.-T. Li. Sampling and summarization for social networks (tutorial), 2013.
- [24] H. Maserrat and J. Pei. Neighbor query friendly compression of social networks. In *SIGKDD*, 2010.
- [25] T. Milo and D. Suciu. Index structures for path expressions. In *ICDT*, 1999.
- [26] S. Navlakha, R. Rastogi, and N. Shrivastava. Graph summarization with bounded error. In *SIGMOD*, 2008.
- [27] M. Palmonari, A. Rula, R. Porrini, A. Maurino, B. Spahiu, and V. Ferme. ABSTAT: linked data summaries with abstraction and statistics. In *ESWC Workshops*, 2015.
- [28] F. Picalausa, Y. Luo, G. H. L. Fletcher, J. Hidders, and S. Vansummeren. A structural approach to indexing triples. In *ESWC*, 2012.
- [29] M. Rudolf, M. Paradies, C. Bornhövd, and W. Lehner. SynopSys: large graph analytics in the SAP HANA database through summarization. In *GRADES*, 2013.

-
- [30] A. Schätzle, A. Neu, G. Lausen, and M. Przyjaciół-Zablocki. Large-scale bisimulation of RDF graphs. In *SWIM*, 2013.
 - [31] B. Spahiu, R. Porrini, M. Palmonari, A. Rula, and A. Maurino. ABSTAT: ontology-driven linked data summaries with pattern minimalization. In *SumPre*, 2016.
 - [32] Y. Tian, R. A. Hankins, and J. M. Patel. Efficient aggregation for graph summarization. In *SIGMOD*, 2008.
 - [33] T. Tran, G. Ladwig, and S. Rudolph. Managing structured and semistructured RDF data using structure indexes. *IEEE TKDE*, 25(9), 2013.
 - [34] T. Tran, H. Wang, S. Rudolph, and P. Cimiano. Top-k exploration of query candidates for efficient keyword search on graph-shaped (RDF) data. In *ICDE*, 2009.
 - [35] G. Troullinou, H. Kondylakis, E. Daskalaki, and D. Plexousakis. RDF digest: Efficient summarization of RDF/S KBs. In *ESWC*, 2015.
 - [36] O. Udrea, A. Pugliese, and V. S. Subrahmanian. GRIN: A graph based RDF index. In *AAAI*, 2007.
 - [37] W3C. Resource description framework. <http://www.w3.org/RDF/>.
 - [38] <http://gromgull.net/blog/2010/09/btc2010-basic-stats>, 2010.
 - [39] X. Yan, P. S. Yu, and J. Han. Graph indexing: A frequent structure-based approach. In *SIGMOD*, 2004.
 - [40] P. Zhao, J. X. Yu, and P. S. Yu. Graph indexing: Tree + delta \geq graph. In *VLDB*, 2007.



**RESEARCH CENTRE
SACLAY – ÎLE-DE-FRANCE**

1 rue Honoré d'Estienne d'Orves
Bâtiment Alan Turing
Campus de l'École Polytechnique
91120 Palaiseau

Publisher
Inria
Domaine de Voluceau - Rocquencourt
BP 105 - 78153 Le Chesnay Cedex
inria.fr

ISSN 0249-6399