



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. 2016. hal-01325900v1

HAL Id: hal-01325900

<https://inria.hal.science/hal-01325900v1>

Submitted on 2 Jun 2016 (v1), 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 — June 2016 — 37 pages

Abstract: The Resource Description Framework (RDF) is the W3C's graph data model for Semantic Web applications. We study the problem of RDF graph summarization: given an input RDF graph G , find an RDF graph H_G which summarizes G as accurately as possible, while being possibly orders of magnitude smaller than the original graph. Summaries are aimed as a help for RDF graph exploration, as well as query formulation and optimization.

We devise four kinds of RDF graph summaries obtained as quotient graphs, with equivalence relations reflecting the similarity between nodes w.r.t. their types or connections. We also study whether they enjoy the formal properties of representativeness (H_G should represent as much information about G as possible) and accuracy (H_G should avoid, to the possible extent, reflecting information that is not in G). Finally, we report the experiments we made on several synthetic and real-life RDF graphs.

Key-words: RDF, data summary

**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. Nous étudions le problème du résumé de graphes RDF : étant donné un graphe RDF G , trouver un graphe RDF H_G résumant G aussi précisément que possible, tout en étant si possible plusieurs ordres de magnitude plus petit que le graphe original. Nos résumés sont destinés à aider l'exploration de graphes RDF, ainsi que la formulation et l'optimisation de requêtes.

Nous proposons quatre sortes de résumé de graphe RDF, obtenus comme des quotients de graphes dont les relations d'équivalence reflètent la similarité entre noeuds vis-à-vis de leurs types ou connexions. Nous étudions aussi s'ils possèdent les propriétés formelles de représentativité (H_G devrait représenter autant d'information de G que possible) et de précision (H_G devrait éviter, autant que possible, de refléter des informations qui ne sont pas dans G). Enfin, nous présentons des expériences faites sur plusieurs graphes RDF synthétiques ou issus d'applications réelles.

Mots-clés : RDF, résumé 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 often *large* and *varied*, produced in a variety of contexts, e.g., scientific applications, social or online media, government data etc. They are *heterogeneous*, i.e., resources described in an RDF graph may have very different sets of properties. An RDF resource may have: no types, one or several types (which may or may not be related to each other). *RDF Schema* (RDFS) information may optionally be attached to an RDF graph, to enhance the description of its resources. Such statements also entail that in an RDF graph, some data is **implicit**. According to the W3C RDF and SPARQL specification, **the semantics of an RDF graph comprises both its explicit and implicit data**; in particular, SPARQL query answers must be computed *reflecting both the explicit and implicit data*. These features make RDF graphs complex, both structurally and conceptually. It is intrinsically hard to get familiar with a new RDF dataset, especially if an RDF schema is sparse or not available at all.

In this work, we study the problem of *RDF summarization*, that is: given an input RDF graph G , find an RDF graph H_G which *summarizes G as accurately as possible, while being possibly orders of magnitude smaller* than the original graph. Such a summary can be used in a variety of contexts: to help an RDF application designer get acquainted with a new dataset, as a first-level user interface, or as a support for query optimization as typically used in semi-structured graph data management [7] etc. Our approach is *query-oriented*, i.e., a summary should enable static analysis and help formulating and optimizing queries; for instance, querying a summary of a graph should reflect whether the query has some answers against this graph, or finding a simpler way to formulate the query etc. Ours is the first semi-structured data summarization approach focused on *partially explicit, partially implicit* RDF graphs.

In the sequel, Section 2 recalls the basics of the RDF data, schema and queries, and sets the requirements for our query-oriented RDF summaries. We then introduce a general concept of RDF summary in Section 3, then decline it into several distinct brands, in Section 4 and 5. We formally establish the properties of these summaries with respect to the representation of the graph they derive from, and analyze the complexity of building them; in all cases, this complexity is in the low-degree polynomials in the size of the graph. We have fully implemented these summarization procedures; we describe our implementation and report experimental results, before discussing related works. Finally, as a picture is worth a thousand words, graphical representations of sample summaries can be found at:

<https://team.inria.fr/cedar/projects/rdfsummary>.

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 $s p o$, stating that the *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 W3C's RDF specification, using uniform resource identifiers (URIs), typed or un-typed literals (constants), and *blank nodes* (unknown URIs or literals) corresponding to a form of incomplete information; they can be seen as some *unknown URI or literal tokens*.

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	$s \leftrightarrow_d o$	$\Pi_{\text{domain}(s)} \subseteq o$
Range typing	$s \leftrightarrow_r o$	$\Pi_{\text{range}(s)} \subseteq o$

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

Notations. We use s , p , and o in triples as placeholders. Literals are shown as strings between quotes, e.g., “string”.

Figure 1 (top) shows how to use triples to describe resources, that is, to describe the type of a resource (unary relation) and a resource property (binary relation). The RDF standard provides 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 RDF graph G shown below describes a book, identified by `doi1`: its author (a blank node `_:b1` related to the author name), title and date of publication.

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

RDF Schema allows enhancing the descriptions in RDF graphs by means of *RDFS triples*, declaring *semantic constraints* between the classes and the properties used in those graphs. Figure 1 (bottom) shows the allowed constraints and how to express them; *domain* and *range* denote respectively the first and second attribute of every property.

This leads to **implicit triples** which are part of the RDF graph even though they are not explicitly present in it. An implicit triple can be obtained by an *immediate entailment step* based on (i) an RDFS constraint (of one of the four forms at the bottom of Figure 1), and (ii) either a second constraint (also called *schema triple*) or an RDF triple that is not a constraint (also termed *data triple*). A triple is *entailed by a graph* G , if and only if there is a sequence of applications of entailment rules that leads from the graph to the triple (where at each step of the entailment sequence, the triples previously entailed are also taken into account). For instance, assume that the RDF graph G above is extended with the following constraints.

- books are publications:
`Book` \prec_{sc} `Publication`
- writing something means being an author:
`writtenBy` \prec_{sp} `hasAuthor`
- `writtenBy` is a relation between books and people:
`writtenBy` \leftrightarrow_d `Book` and
`writtenBy` \leftrightarrow_r `Person`

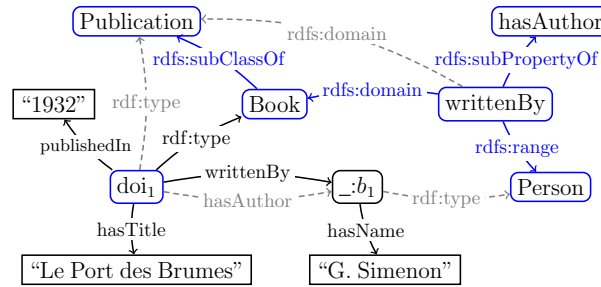


Figure 2: Sample RDF graph.

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

Well-behaved RDF graphs. We say an RDF graph is *well-behaved* if (i) no class appears in the property position; (ii) no class has properties other than the `rdf:type` and than those of RDF Schema. These constraints are not part of the RDF specification, yet they reflect common-sense data design decisions, and hold in all datasets we experimented with. In the sequel, we assume all graphs are all-behaved.

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

The triple-based representation of an RDF graph. We see an RDF graph G as a *disjoint union of three components* $G = \langle D_G, S_G, T_G \rangle$

the schema S_G , that is the set of all G triples whose properties are \prec_{sc} , \prec_{sp} , \leftrightarrow_d or \leftrightarrow_r ;

the type T_G , that is the set of the τ triples from G ;

the data D_G , which holds all the remaining triples of G .

Note that each of D_G , S_G , and T_G is an RDF graph by itself.

The graph-based representation of an RDF graph. As per the RDF specification [19], *the set of nodes* of an RDF graph is the set of subjects and objects of triples in the graph. Further, we define: (i) a *data node* as any URI or literal occurring as a subject or object in D_G , or as a subject in T_G ; (ii) a *class node* as a URI appearing in the object position of triples from T_G , and (iii) a *property node* as a URI appearing in the subject or object position of \prec_{sp} triples, or in the subject position of \leftrightarrow_d or \leftrightarrow_r triples from S_G . The *edges* of an RDF graph correspond to its triples, where the subject/object is the edge source/target node, and the property is the edge label. More specifically, we differentiate between *data*, *type* and *schema edges*, based on their labeling property.

Size and cardinality notations. We denote by $|G|_n$ the number of nodes in a graph G , and by $|G|_e$ 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 .

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 fixed-point of repeatedly applying these 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). Clearly, a graph G entails a triple if and only if

the triple belongs to G^∞ . RDF entailment is part of the RDF standard; *the answers to a query posed on G must take into account all triples in G^∞ , since the semantics of an RDF graph is its saturation.*

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 [12]. 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 the variables 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 G 's explicit triples, thus may lead to an incomplete answer; the complete answer is obtained by evaluating q against G^∞ . For instance, the query below asks for name of the author of “Le Pont des Brumes”:

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

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.

2.2 RDF Summary Requirements

We assume that *the summary H_G of an RDF graph G is an RDF graph itself*. Further, we require summaries to satisfy the following conditions:

Schema independence It must be possible to summarize G whether or not it has a schema.

Semantic completeness The summary of G must reflect not only the explicit data, but also the implicit one, 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 orders of magnitude.

Representativeness The summary should preserve as much information about G as possible¹.

Accuracy The summary should avoid, to the extent possible, reflecting data that does not exist in G .

Criteria for representativeness and accuracy. Our query-oriented RDF graph summarization leads us to the following criteria. For *representativeness*, queries with results on G should also have results on the summary. Symmetrically, for *accuracy*, a query that can be matched on the summary, should also be matched on the RDF graph itself.

To formalize our criteria, we use \mathcal{Q} to denote an *RDF query language (dialect)*; a concrete choice of such a dialect will shortly follow.

¹A clear trade-off exists between compactness and representativeness; we present concrete proposals for such a trade-off later on.

Definition 1 (QUERY-BASED REPRESENTATIVENESS) *Let G be any RDF graph. H_G is \mathcal{Q} -representative of G if and only if for any query $q \in \mathcal{Q}$ such that $q(G^\infty) \neq \emptyset$, we have $q(H_G^\infty) \neq \emptyset$.*

Note that several graphs may have the same summary, which corresponds to the intuition that a summary loses some of the information from the original graph. If two RDF graphs differ only with respect to such information, they have the same summary. We term *inverse set* of a summary H_G , the set of all RDF graphs whose summary is H_G . This leads to the accuracy criterion:

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

For the sake of compactness, we decide that the (voluminous) set of literals, along with subject and object URIs for non- τ triples from G should not appear in H_G . However, given that property URIs are often specified in SPARQL queries [1], and that there are typically far less distinct property URIs than there are distinct subject or object URIs [20], property URIs should be preserved by the summary. This leads us to the following SPARQL dialect:

Definition 3 (RELATIONAL BGP QUERY) *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.*

A sample RBGP is:

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

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 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, one can see that 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 (R)BGP 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.)

Let RBGP* denote an extension of the RBGP dialect whereas in property positions, one can either use URIs or variables. We find that:

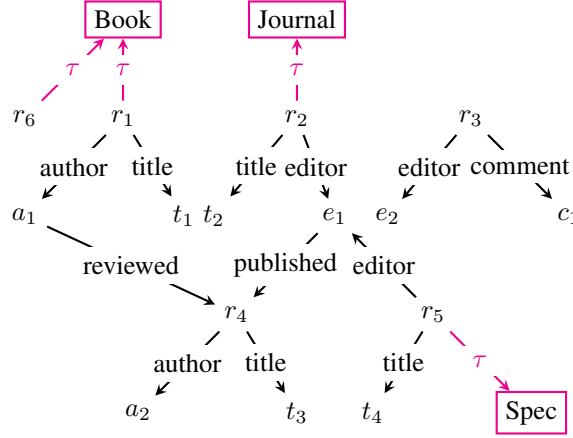


Figure 3: Sample RDF graph.

Proposition 1 *RBGP representativeness entails RBGP* representativeness.*

Proof 1 Let q^* be an RBGP* query which is non-empty on G^∞ and H_G be an RBGP-representative summary of G . We show that non-emptiness of q^* on G^∞ entails its non-emptiness on H_G^∞ .

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 H_G is RBGP representative, q also has results on H_G^∞ . Therefore, $q(H_G^\infty) \neq \emptyset$, and given that $q \subseteq q^*$ (query containment), it follows that $q^*(H_G^\infty) \neq \emptyset$.

Proposition 2 *RBGP accuracy entails RBGP* accuracy.*

Proof 2 Let H_G be an RBGP-accurate summary of G and q^* an RBGP* query which is non-empty on H_G^∞ . We show that non-emptiness of q^* on H_G^∞ entails its non-emptiness on G'^∞ , where G' belongs to the inverse set of H_G .

Given that $q^*(H_G^\infty)$ is non-empty, there exists at least an embedding of q^* into H_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 H_G^∞ . Clearly, q has results on H_G^∞ , and since H_G is RBGP accurate, q also has results on G'^∞ , where G' belongs to the inverse set of H_G . Therefore, $q(G'^\infty) \neq \emptyset$, and given that $q \subseteq q^*$ (query containment), it follows that $q^*(G'^\infty) \neq \emptyset$.

Based on these properties, we shall work with RBGPs to establish properties also for RBGP*.

3 RDF summarization

Let $G = \langle D_G, S_G, T_G \rangle$ be an RDF graph. For illustration, we will rely on the graph in Figure 3; here and in the sequel, we show class nodes in purple boxes, while T_G triples appear as purple arrows.

We recall a classical definition from graph theory:

Definition 4 (QUOTIENT GRAPH) Let $G = (V, E, \lambda)$ be a labeled directed graph whose vertices are V , whose edges are $E \subseteq V \times V$ and whose edges are labeled by a function $\lambda : E \rightarrow A$ for a given label set A . Let $\equiv \subseteq V \times V$ be an equivalence relation over the nodes of V .

The quotient graph of G over \equiv , denoted G_{\equiv} , is a 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 $n_1 \xrightarrow{a} n_2 \in E$.

We call *data property* any property p occurring in D_G , and *data triple* any triple in D_G . For instance, in Figure 3, the data properties are: author, title, editor, comment, reviewed, and published, which for the sake of brevity we will denote as a , t , e , c , r and p .

Summarization approach and node equivalence. We will define summaries through quotient graphs, using various equivalence relations. To establish such relations, we first examine the data properties of graph nodes; we identify interesting relations between these properties (Section 3.1) and use these to define *node equivalence relations based on their data properties* in Section 3.2, where we also introduce a simple *node equivalence relation based on node types*. In Section 3.3, we formalize RDF summaries, and study some of their important properties.

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 5 (PROPERTY RELATIONS AND CLIQUES) *Let p_1, p_2 be two G data properties, and r_1, r_2 be two data nodes in G .*

1. $p_1, p_2 \in G$ are source-related iff one of the following conditions holds: (i) a resource in G has both p_1 and p_2 , or (ii) G holds a resource r and a data property p_3 such that r has p_1 and p_3 , and moreover p_3 and p_2 are source-related.
2. $p_1, p_2 \in G$ are target-related in the same conditions as above (replacing in (i) “has a property” by “is the value of a property” and in (ii) “source” by “target”).
3. 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 ; this entails also that a and e are source-related. Properties r and p are target-related due to r_4 . The source cliques are:

$$SC_1 = \{a, t, e, c\}; SC_2 = \{r\}; SC_3 = \{p\};$$

whereas the target cliques are:

$$TC_1 = \{a\}; TC_2 = \{t\}; TC_3 = \{e\}; TC_4 = \{c\}; TC_5 = \{r, p\}$$

A source clique can be viewed as “all the data properties of resources of the same kind”; in the above example, 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 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.

Definition 6 (PROPERTY DISTANCE IN A CLIQUE) *The distance between two data properties p, p' in a source clique SC is:*

- 0 if there exists a resource in G having them both;

r	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
	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
	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.

- otherwise, the smallest integer n such that G holds resources $r_0, r_2, \dots, r_n \in G$ and data properties p_1, \dots, p_n such that r_0 has p and p_1, r_1 has p_1 and p_2, \dots, r_n has p_n and p' .

For instance, in our sample graph, the distance between a and t is 0 since r_1 has them both. The distance between a and e is 1, while the distance between a and c is 2.

It is easy to see that p and p' are at distance n for $n > 0$ iff there exists a resource having p and p' , and further p'' is at distance $n - 1$ from p' .

Source and target cliques are defined w.r.t. a given graph G ; when moving from G to G^∞ , resources may have more data properties due to the \prec_{sp} constraint, thus possibly different cliques. The following important property characterizes the relationship between the cliques of G and G^∞ :

Lemma 1 (Saturation vs. property cliques) *Let C, C_1, C_2 be distinct source (or target) cliques corresponding to G .*

1. There exists exactly one source (resp. target) clique C^∞ corresponding to G^∞ such that $C \subseteq C^\infty$.
2. We call saturated clique and denote C^+ the set of properties comprising all the properties in C and all their generalizations (superproperties). 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. Let p_1, p_2 be two data properties in different source (or target) cliques of G , namely 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 source cliques of G , $k \geq 0$, denoted D_1, D_2, \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 3 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.

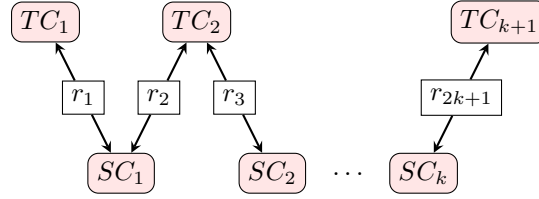


Figure 4: Alternating sequence of shared cliques between weakly related resources $r_1, r_2, \dots, r_{2k+1}$.

3. “Only if”: We make the proof by induction on the distance between p_1 and p_2 in G^∞ .

If the distance is 0, they belong to the same resource in G^∞ , but clearly they cannot belong to the same resource in G given that their cliques are distinct there. Then, it must be that for a property p , the schema of G implies that $p_1 \prec_{sp} p$ and $p_2 \prec_{sp} p$, and further in G , a resource r has both p_1 and p_2 . This entails that C_1^+ and C_2^+ intersect on p .

Now assume $k > 0$, and p_1 and p_2 are at distance k in G^∞ . In this case, there exists a property p_3 and a resource r such that r has p_1 and p_3 in G^∞ , whereas p_2 and p_3 are at a distance $n < k$ in G^∞ . Thus, by the induction hypothesis, there exist n source cliques D_1, D_2, \dots, D_n such that:

$$C_2^+ \cap D_1^+ \neq \emptyset, D_1^+ \cap D_2^+ \neq \emptyset, \dots, D_n^+ \cap C_3^+ \neq \emptyset$$

where C_3 is the source clique of p_3 . Since r has p_1 and p_3 in G^∞ , it must be the case that r holds two properties p'_1, p'_3 in G such that after (possibly) generalization, i.e., applying \prec_{sp} , one leads to p_1 while the other leads to p_3 . Hence, the source clique C_r of r in G comprises both p'_1 and p'_3 , so C_r^+ contains p'_1, p'_3, p_1 and p_3 . Therefore, C_r^+ intersects C_1^+ on p_1 , establishing a chain of $k + 1$ saturated cliques such intersecting at one end C_1^+ and at the other end C_2^+ . This finalizes our proof by induction.

“If”: Applying item 2. of the lemma for C_1^+ and D_1^+ , if these intersect, then their properties are in the same G^∞ clique, and the same holds for each pair of successive saturated cliques; thus, all the properties in $C_1, D_1, \dots, D_k, C_2$ are in the same G^∞ clique. Thus, p_1 and p_2 are in the same G^∞ clique.

3.2 Equivalence relations among graph nodes

We use the above source and target cliques to define two notions of equivalence between data nodes of a graph:

Definition 7 (STRONG AND WEAK EQUIVALENCE) *Two data nodes of G are strongly related, denoted $n_1 \equiv_s n_2$, iff they have the same source and target cliques.*

Two nodes are weakly related, denoted $n_1 \equiv_w n_2$, iff they have either a same non-empty source or target cliques, or they are both related to another data node of G .

Observe that strong relatedness entails weak relatedness.

In Figure 3, the resources r_1, r_2, r_3, r_5 are *strongly* related to each other, as well as t_1, t_2, t_3, t_4 . Moreover, r_1 to r_5 are *weakly* related 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, resources can also be weakly related through an *alternating sequence of source and target cliques*, as illustrated in Figure 4.

If one does not wish to consider data properties as a basis for node equivalence, a simple equivalence based on node types is:

Definition 8 (TYPE-BASED EQUIVALENCE) *Two nodes n_1, n_2 of G are type-equivalent, denoted $n_1 \equiv_T n_2$, iff (i) both n_1 and n_2 have types in G and (ii) they have exactly the same set of types.*

Recall that some or even all nodes in an RDF graph may lack types; type-based equivalence only Clearly, each equivalence relation defines a partition over the data nodes in G .

3.3 RDF summaries

We define RDF summaries based on graph quotients:

Definition 9 (RDF SUMMARY) *Given an RDF graph $G = \langle D_G, S_G, T_G \rangle$ and an equivalence relation \equiv among G 's nodes, a summary of G is an RDF graph $H_G = \langle D_H, S_H, T_H \rangle$ whose triples are defined as follows:*

SCH H_G has the same schema triples as G : $S_H = S_G$;

TYP+DAT $T_H \cup D_H$ is the quotient of $T_G \cup D_G$ by \equiv : T_H is the set of τ triples in $(T_H \cup D_H)_{\equiv}$, while D_H holds the remaining triples.

This definition *preserves the RDF schema unchanged*. We do this, first, because a schema (when present) provides a valuable and typically compact set of constraints describing the data semantics; and second, because we want the summary to enable, on a smaller graph, the kinds of reasoning possible on the original graph – in particular, as we shall see, as a means to achieve summary completeness (Section 2.2). Further, since the summary is an RDF graph, *all its nodes must be labeled by URIs or literals*; in particular, the subjects in T_H and D_H must be URIs.

We find that:

Proposition 3 (SUMMARY REPRESENTATIVENESS) *An RDF summary H is RBGP-representative.*

Proof 4 *Let q be a query such that $q(G^\infty) \neq \emptyset$; we need to show that $q(H_G^\infty) \neq \emptyset$.*

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 $(H_G)^\infty$.

First, consider a triple pattern t of q whose embedding $\phi(t) \in G^\infty$ also belongs to G . If $\phi(t) = s \ p \ o \in G$, that is, $\phi(t)$ is a data triple, H_G holds the triple $f(s) \ p \ f(o)$, thus $(H_G)^\infty$ also comprises it. If $\phi(t)$ was a type triple of the form $s \ \tau \ c$, the triple $f(s) \ \tau \ c$ belongs to H_G , thus also to $(H_G)^\infty$.

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

First, $\phi(t')$ may be data triple, and let $t_d = s_d \ p_d \ o_d$ be a data triple in G and t_s be a schema triple in G^∞ such that $\phi(t')$ is directly entailed by t_d and t_s . As explained above, $f(s_d) \ p_d \ f(o_d) \in H_G$; at the same time, t_s also belongs to $(H_G)^\infty$. Given that the latter includes all the triples that it entails, 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 $(H_G)^\infty$.

The second case is when $\phi(t')$ is a τ triple. 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 \hookrightarrow_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 $t_s = c_1 \ \prec_{sc} \ c_2 \in (S_G)^\infty$. In this case, H_G holds the triple $f(s) \ \tau \ c_1$ which is also present in $(H_G)^\infty$, thus the same inference step applies in $(H_G)^\infty$ to produce $f(s) \ \tau \ c_2$.*

Thus, a q triple (whether its property is a data property, or τ) 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 $(H_G)^\infty$, concluding our proof.

We are interested in summaries having the following property:

Definition 10 (FIXPOINT PROPERTY) *A summary H has the fixpoint property iff for any graph G , $H_{H_G} = H_G$ holds.*

Intuitively, the fixpoint property expresses the fact that a summary cannot be summarized further, i.e., H_G is its own summary. This is a desirable property as we wish our summaries to be as compact as possible, while satisfying our requirements. It turns out that our summaries enjoy this property, as they are quotient graphs:

Proposition 4 (SUMMARY FIXPOINT) *An RDF summary has the fixpoint property.*

Proof 5 *Suppose that a summary does not have the fixpoint property, i.e., $H_{H_G} \neq H_G$. Observe first that, by definition of quotient graph, if H_{H_G} and H_G differ, then H_{H_G} must have less nodes than H_G , since the edges are completely determined by the set of nodes.*

Case \equiv_S : if two H_G 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, every 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 end up in a single H_G node, a contradiction.

Case \equiv_W : if H_{H_G} has less nodes than H_G , then there exist two H_G nodes n_1, n_2 that have the same source or target cliques. In this case, there exists some property common to these same cliques, hence all the G nodes represented n_1 and n_2 end up in a single H_G node by definition of a graph quotient using \equiv_W , a contradiction.

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

Moreover, the fixpoint property of our summaries has another immediate good consequence, as a summary is obviously a summary of itself:

Proposition 5 (ACCURACY) *An RDF summary is accurate.*

Proof 6 *The proof follows from the fact that any summary has a fixpoint property: a summary H which is its own summary, corresponds at least to the RDF graph H itself.*

Finally, the size of our summaries is bounded by that of the summarized graphs:

Proposition 6 (SUMMARY SIZE BOUND) *The size of an RDF summary of an RDF graph G is at most the size of G .*

Proof 7 *This bound is reached, i.e., $G = H_G$ holds, when:*

\equiv_W **case:** *all nodes have distinct source and target properties,*

\equiv_S **case:** *all nodes have distinct source or target cliques,*

\equiv_T **case:** *all nodes have distinct type sets.*

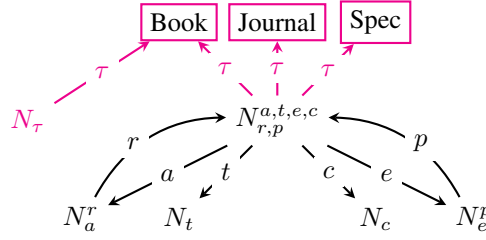


Figure 5: Weak summary of the RDF graph in Figure 3.

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 simplest summary is solely based on weak equivalence:

Definition 11 (WEAK SUMMARY) *The weak summary of the graph G , denoted W_G , is its quotient graph w.r.t. the weak equivalence relation \equiv_w .*

It follows from the definition of quotient graphs that for each G node equivalence set w.r.t. \equiv_w , there is exactly one node in W_G . Further, note that the partition of G 's nodes into sets according to \equiv_w is also a *partition of G data properties at the source*, that is: in G , the sources of 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 weakly equivalence class $\{r_1, r_2, r_3, r_4, r_5\}$. By the same reasoning, the \equiv_w partition over G 's node induces a *partition of G 's source cliques*; and by a symmetrical reasoning, it introduces also a *partition over G 's target cliques*. 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 nodes in G one can associate through a bijection, a set of G source cliques, and a set of G target cliques. We shall refer to the node $n_S \in W_G$ representing S as:

$$n_S = N\left(\bigcup_{n \in S} TC(n), \bigcup_{n \in S} SC(n)\right)$$

where N , termed *representation function*, is any injective function taking as input two sets of URIs (respectively, a set of target data properties, and a set of source data properties), and returning a new URI. Without loss of generality, we will use N to denote a concrete representation function which assigns URIs to nodes in quotient (RDF) graphs.

Notations. Whenever this simplifies reading, we may use N_r or N_{SC}^{TC} (where TC is the first and SC is the second input to N), to denote $f_w(r)$. Further, for simplicity, we will mostly *omit the set delimiters* when showing TC and SC , and omit one such set altogether if it is empty.

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 5. 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 $\{t_1, t_2, t_3, t_4\}$;
- N_e^p for the set $\{e_1, e_2\}$;
- N_c for the set $\{c_1\}$.

The edges from $N_{r,p}^{a,t,e,c}$ to N_a and N_t are due to r_1 and item **DAT** of the summary definition (Definition 9); 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 7 (UNIQUE DATA PROPERTIES) *Each \mathbb{G} data property appears exactly once in $\mathbb{W}_{\mathbb{G}}$.*

Proof 8 *First, note that the sets of properties of which weak nodes are targets, are pairwise disjoint. Indeed, if this was not the case, one target clique TC would have needed to be unioned in two distinct weak summary nodes; this supposes the existence of two \mathbb{G} nodes r_1, r_2 which have the same target clique but are part of two distinct relatedness partition sets, which is not possible. The same holds for the sets of properties of which weak nodes are sources. Thus, any data property has at most one source and at most one target in $\mathbb{W}_{\mathbb{G}}$. As a consequence, and due to rule **DAT** stating how summary edges are built, there is exactly one p -labeled edge in $\mathbb{W}_{\mathbb{G}}$ for every data property in \mathbb{G} .*

Importantly, the above Property 7 entails that the number of data edges in $\mathbb{W}_{\mathbb{G}}$ is exactly $|\mathbb{D}_{\mathbb{G}}|_p^0$, the number of distinct data properties in \mathbb{G} . Thus, its number of data nodes is at most $2|\mathbb{D}_{\mathbb{G}}|_p^0$. The number of type triples in $\mathbb{W}_{\mathbb{G}}$ is bound by $\min(|\mathbb{T}_{\mathbb{G}}|_e, 2|\mathbb{D}_{\mathbb{G}}|_p^0 * |\mathbb{T}_{\mathbb{G}}|_o^0)$: the latter corresponds to the case when every data node in $\mathbb{W}_{\mathbb{G}}$ is of every type in $\mathbb{T}_{\mathbb{G}}$.

The next important property of the weak summary is:

Proposition 8 (WEAK COMPLETENESS) *For a given graph \mathbb{G} , $\mathbb{W}_{\mathbb{G}^\infty} = \mathbb{W}_{(\mathbb{W}_{\mathbb{G}})^\infty}$ holds.*

This property is important, as it gives a mean to compute $\mathbb{W}_{\mathbb{G}^\infty}$ without saturating \mathbb{G} , but only summarizing \mathbb{G} , then saturating the smaller (typically by several orders of magnitude) $\mathbb{W}_{\mathbb{G}}$, and applying again weak summarization. This is exemplified in Figure 6, which traces the transformation of a graph \mathbb{G} on one hand, into $\mathbb{W}_{\mathbb{G}}$, $(\mathbb{W}_{\mathbb{G}})^\infty$, and then $\mathbb{W}_{(\mathbb{W}_{\mathbb{G}})^\infty}$; and on the other hand, from \mathbb{G} to \mathbb{G}^∞ then $\mathbb{W}_{\mathbb{G}^\infty}$. The graph at the bottom right in the figure illustrates the equality stated by the proposition.

The proof of Proposition 8 requires the following crucial lemma:

Lemma 2 (Shared cliques vs. summarization) *If two resources r', r'' share a target clique in \mathbb{G}^∞ , their nodes $N_{r'}, N_{r''}$ representing them, share a target clique in $(\mathbb{W}_{\mathbb{G}})^\infty$.*

Proof 9 *Let TC', TC'' be the target cliques of r', r'' in \mathbb{G} . It follows from Lemma 1 (“only if”) that \mathbb{G} holds a set of target cliques D_1, \dots, D_k such that*

$$TC'^+ \cap D_1^+ \neq \emptyset, D_1^+ \cap D_2^+ \neq \emptyset, \dots, D_{k-1}^+ \cap D_k^+ \neq \emptyset, D_k^+ \cap TC''^+ \neq \emptyset \quad (*)$$

Let $N_{r'}, N_{r''}$ be the nodes representing r', r'' in $\mathbb{W}_{\mathbb{G}}$. We know that:

$$\begin{aligned} TC' &= TC(r', \mathbb{G}) \subseteq TC(N_{r'}, \mathbb{W}_{\mathbb{G}}) \\ TC'' &= TC(r'', \mathbb{G}) \subseteq TC(N_{r''}, \mathbb{W}_{\mathbb{G}}) \end{aligned}$$

These inequalities above also hold after saturating the cliques:

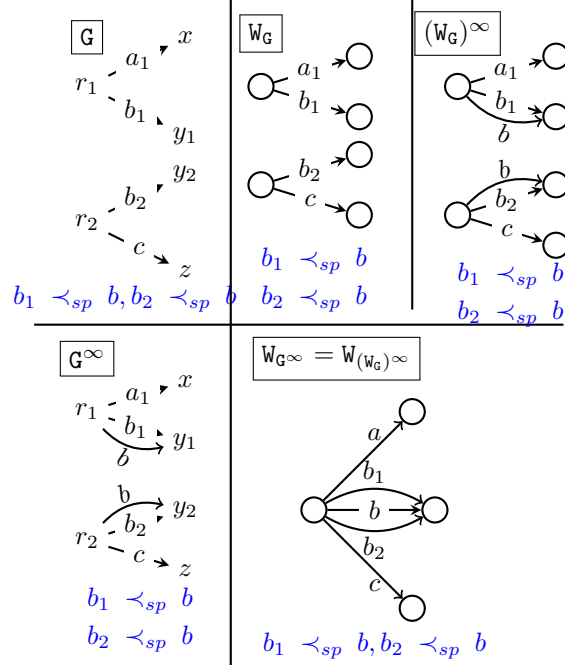


Figure 6: Weak summary completeness illustration (summary nodes shown as unlabeled circles).

$$\begin{aligned}
 TC'^+ &= TC(r', G)^+ \subseteq TC(N_{r'}, W_G)^+ \\
 TC''^+ &= TC(r'', G)^+ \subseteq TC(N_{r''}, W_G)^+
 \end{aligned}$$

Thus, if a sequence of saturated target cliques of G intersecting each other connect TC'^+ to TC''^+ (established in (*) above), a similar sequence of saturated target cliques of W_G connects $TC(N_{r'}, W_G)^+$ to $TC(N_{r''}, W_G)^+$ in W_G . Thus, applying again Lemma 1 ("if"), these lead to the same target clique TC in $(W_G)^\infty$.

Thus, $N_{r'}$ and $N_{r''}$ have the same target clique in $(W_G)^\infty$.

Based on the lemma, we can now show:

Proof 10 (Proposition 8) The schema components are obviously identical, and equal to $(S_G)^\infty$. Below we focus on the data components.

To prove that the weak summaries of G^∞ and $(W_G)^\infty$ have the same data components, we start by showing that the nodes of the data components are obtained by identical calls to N . This reduces to showing the following two statements:

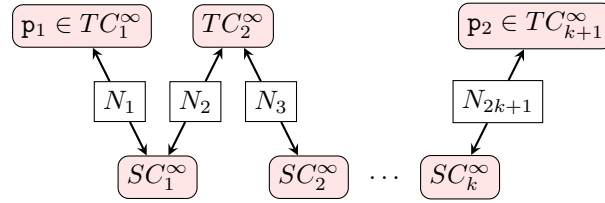
1. two G properties appear in the target set of a given node in W_{G^∞} iff they appear in the target set of a given node in $W_{(W_G)^\infty}$ (and similarly for the source sets);
2. a target and a source clique of G are in the same data node in W_{G^∞} iff they are in the same data node in $W_{(W_G)^\infty}$.

Together, 1. and 2. above ensure that the data properties of G are "grouped together identically" in source, respectively target sets, when building the nodes of the two weak summaries.

We show 1. for target properties. By the weak summary definition, two data properties appear in the target set of the same node in \mathbb{W}_G iff they are target-related. Thus, 1. reduces to showing that: p_1, p_2 are target-related in G^∞ iff they are target-related in $(\mathbb{W}_G)^\infty$.

“If”: Two properties p_1, p_2 are target-related in $(\mathbb{W}_G)^\infty$ iff $(\mathbb{W}_G)^\infty$ holds a set of nodes $N_1^\infty, \dots, N_{2k+1}^\infty$ and some target and source cliques $TC_1^\infty, \dots, TC_{k+1}^\infty, SC_1^\infty, \dots, SC_k^\infty$ such that:

$p_1 \in TC_1^\infty$, N_1 has the target clique TC_1^∞ and the source clique SC_1^∞ , N_2 has the source clique SC_1^∞ and the source clique TC_2^∞ , N_3 has the target clique TC_2^∞ and the source clique SC_2^∞ , and so on until N_{2k+1} has the target clique TC_{k+1}^∞ and the source clique SC_k^∞ , $p_2 \in TC_{k+1}^\infty$ (see the sketch below).



We first show that:

- (i) p_1 is target-related in G^∞ to any data property from TC_1^∞ .

To see why, note that TC_1^∞ is the union of the saturations of a set of (disjoint) \mathbb{W}_G cliques $(TC_1^1)^+, (TC_1^2)^+, \dots, (TC_1^{m_1})^+$. Then:

- For every $1 \leq j \leq m_1$, TC_1^j is the target clique of a \mathbb{W}_G node, created in the weak summary exactly to represent a set of related resources which are target of the properties from TC_1^j .
- Thus, in G^∞ , also, all properties from $(TC_1^j)^+$ are target-related, $1 \leq j \leq m_1$.
- Further, given that each $(TC_1^j)^+$ intersects $(TC_1^{j+1})^+$ for $1 \leq j < m_1$, and that all properties in the TC_1^j are represented in G for $1 \leq j \leq m_1$, it follows that the target properties in all the $(TC_1^j)^+$ (p_1 is among them) are related to each other in G^∞ .

By a similar reasoning:

- (ii) All the properties in a clique TC_i^∞ , with $1 \leq i \leq k+1$, are related between themselves in G^∞ ; in particular, p_2 is related to all the properties in TC_{k+1}^∞ . The same holds for all properties in SC_i^∞ , for $1 \leq i \leq k$.

Next, we show that:

- (iii) The properties in a source clique SC^∞ and a target clique TC^∞ which belong to the same node N_i , $1 \leq i \leq 2k+1$, are data-related in G^∞ .

where two properties are data-related in a graph iff they are sources and/or targets of two related resources from that graph. (This generalizes both source- and target-related properties (Definition 5.) As above, we know that:

$$\begin{aligned} SC^\infty &= (SC^1)^+ \cup (SC^2)^+ \cup \dots \cup (SC^m)^+ \text{ and} \\ TC^\infty &= (TC^1)^+ \cup (TC^2)^+ \cup \dots \cup (TC^n)^+ \end{aligned}$$

for some \mathbb{W}_G source cliques SC^1, \dots, SC^m and target cliques TC^1, \dots, TC^n . We know the data properties in $(SC^1)^+ \cup \dots \cup (SC^m)^+$ are related in \mathbb{G}^∞ (as shown above for $(TC_1^1)^+ \cup \dots \cup (TC_1^{m_1})^+$); similarly those of $(TC^1)^+ \cup \dots \cup (TC^n)^+$ are related in \mathbb{G}^∞ . Moreover, N_i was created in \mathbb{W}_G out of related G nodes having as source clique one among SC^1, \dots, SC^m and as target clique one among TC^1, \dots, TC^n . In \mathbb{G}^∞ , these nodes connect the data properties of SC^∞ with those of TC^∞ .

From (ii) and (iii) above, it follows that all properties in TC_1^∞ are data-related in \mathbb{G}^∞ to those of TC_2^∞ , thus in particular p_1 is data-related to p_2 in \mathbb{G}^∞ .

“Only if”: We assume that p_1, p_2 are target-related in \mathbb{G}^∞ . Thus, \mathbb{G}^∞ features some resources r'_1, \dots, r'_{2k+1} , source cliques $SC_1'^\infty, \dots, SC_k'^\infty$ and target cliques $TC_1'^\infty, \dots, TC_{k+1}'^\infty$ such that:

- r_1 has the target clique TC_1^∞ and the source clique $SC_1'^\infty$; we denote this $r_1(TC_1^\infty, SC_1'^\infty)$ for brevity;
- $r'_1(TC_1'^\infty, SC_1'^\infty), r'_2(TC_1'^\infty, SC_2'^\infty)$ and so on until
- $r'_k(TC_k'^\infty, SC_{k+1}'^\infty),$
- $r_2(TC_k^\infty, SC_{k+1}'^\infty).$

Applying Lemma 2 to each consecutive pair of resources above, i.e., r_1 and r'_1 , then r'_1 and r'_2 etc. we obtain that the \mathbb{W}_G nodes representing each such pair of resources share a clique in \mathbb{W}_G , thus also in $(\mathbb{W}_G)^\infty$. Thus, a similar data-relation path can be exhibited in $(\mathbb{W}_G)^\infty$, relating p_1 to p_2 there.

The above has established that the sets of G target properties given as input to N when creating \mathbb{W}_{G^∞} and $\mathbb{W}_{(\mathbb{W}_G)^\infty}$ are the same, and the same for the source properties.

The item 2. above (guaranteeing the same pairs of G source and target property sets are used on both sides) directly follows from the observation that both summarization and saturation preserve source-target associations, in the following sense. Let ψ denote a transformation that is either summarization or saturation. ψ preserves source-target associations iff: for any node r having a source property p_s and a target property p_t , p_s is also a source property for $\psi(r)$, while p_t is also a target property p_t for $\psi(r)$.

Weak summary nodes completely determine the data edges in the summary (cf. unicity of every data property, Proposition 7). Thus, the two summaries have the same data triples.

Further, for any node $r \in G$, a G triple of the form $r \tau c$ leads to the corresponding triple $\phi_N(r) \tau c$ in \mathbb{W}_G which carries over to $(\mathbb{W}_G)^\infty$ and the corresponding node in $\mathbb{W}_{(\mathbb{W}_G)^\infty}$; similarly, for any inferred triple in \mathbb{G}^∞ , a similar inference will add triples to $(\mathbb{W}_G)^\infty$. Thus, the τ triples in \mathbb{W}_{G^∞} and $\mathbb{W}_{(\mathbb{W}_G)^\infty}$ are the same.

Figure 6 illustrates weak summary completeness on an example; to avoid clutter, summary nodes are shown in the figure as unlabeled circles.

4.2 Typed weak summary

In this section, we introduce a summary based on the weak relatedness notion but which gives a strong preeminence to type informations: it represents together nodes having the same set of RDF types, and quotients the data triples only for nodes having no types.

Formally, given a graph G , we term *typed resources* of G and denote TR_G the set of subjects of T_G (type) triples. We denote *untyped resources* of G the set UN_G of subjects or objects of D_G triples, which have no types. UN_G may comprise URIs and/or literals. The *untyped data graph* of G , denoted UD_G , is the subset of D_G triples whose subject and object belong to UN_G .

We start by introducing two helper notions:

Definition 12 (TYPE-BASED SUMMARY) *The type-based summary of G , denoted T_G , is the summary of G through the \equiv_T equivalence relation.*

This summary 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, given as input a set X of (class) URIs, returns a URI $C(X)$, and given an empty set of URIs, returns a new URI on every call. An easy bijection holds between the \equiv_T equivalence classes of typed resources in G , and their respective sets of types. Thus, every node in the summary T_G can be seen as obtained through a call to $C(X)$, where X is the set of types of the resources in an equivalence class wrt \equiv_T . Figure 7 illustrates a typed summary.

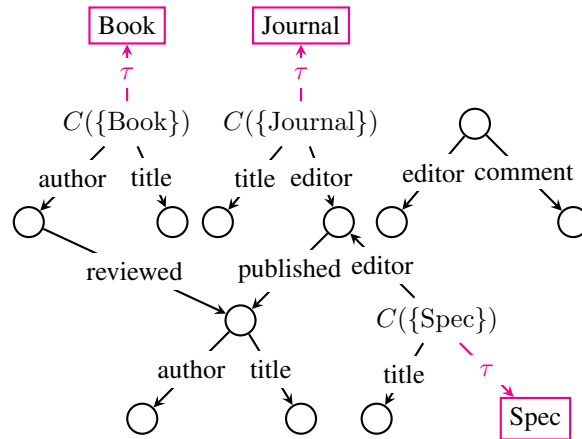


Figure 7: Typed summary of the RDF graph in Figure 3. Each circle is a distinct node whose URI results from calling $C(\emptyset)$.

Definition 13 (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 UW_G , is its summary through \equiv_{UW} .

Untyped weak equivalence restricts weak equivalence to untyped resources only. The untyped-weak summary UW_G summarizes UD_G in weak fashion, and leaves all typed resources untouched. We can now define:

Definition 14 (TYPED WEAK SUMMARY) *The typed weak summary TW_G of an RDF graph G is the untyped-weak summary of the type-based summary of G , namely UW_{T_G} .*

For example, Figure 8 shows the typed weak summary of the RDF graph in Figure 3. Note how distinct types lead to different nodes in the summary.

Proposition 9 (TYPED WEAK FIXPOINT) *The typed weak summary (Definition 14) has the fixpoint property.*

Proof 11 *The proof follows from the fact that both T and UW summarization enjoy the fixpoint property (by Property 4).*

Proposition 10 (TYPED WEAK NON-COMPLETENESS) *There exists an RDF graph G such that $TW_{G^\infty} \neq TW_{(TW_G)^\infty}$ holds.*

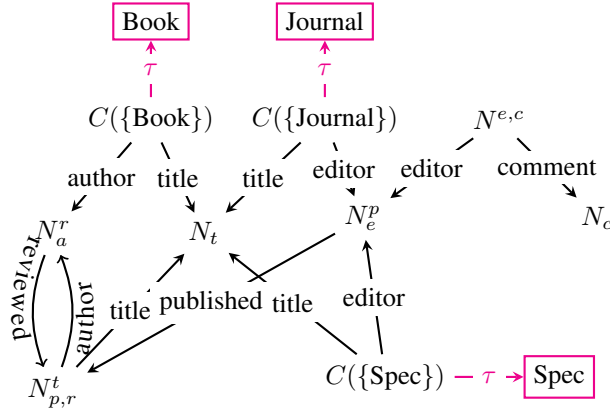
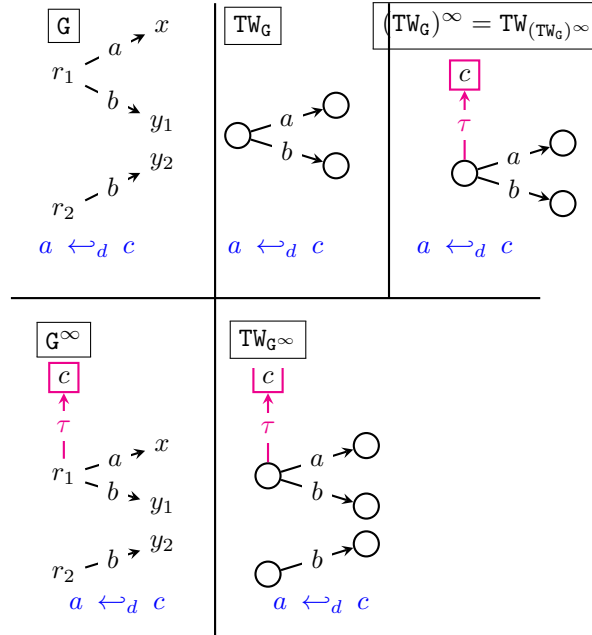
Figure 8: Typed weak summary example: $TW_G = UW_{T_G}$ for the T_G in Figure 7.

Figure 9: Typed weak summary completeness counter-example.

Proof 12 A counter-example appears in Figure 9. In G and TW_G , resources are untyped; a typed resource only appears due to saturation and the schema rule $a \leftrightarrow_d c$. Thus, in TW_G and its saturation, a single node represents all subjects of data properties. In contrast, in G^∞ , one resource is typed and the other one is untyped, leading to distinct nodes in TW_{G^∞} .

More generally, in the presence of domain (\leftrightarrow_d) or range (\leftrightarrow_r) RDF schema rules, one cannot compute TW_{G^∞} from G because the typed weak summary represents typed resources differently from the untyped ones. The \leftrightarrow_d and \leftrightarrow_r schema rules may turn untyped resources into typed ones, thus leading to divergent representations of G 's data nodes in TW_G and respectively TW_{G^∞} .

5 Strong equivalence summaries

In this section, we discuss summary alternatives based on the strong equivalence \equiv_s between graph nodes. Strong summaries (Section 5.1) are mainly based on nodes' data properties; the typed strong summary (Section 5.2) gives preeminence to types.

5.1 Strong summary

RDF data nodes represented by the same strong summary node have *the same source clique and the same target clique*. Formally, based on the N function introduced in Section 4.1, we define:

Definition 15 (STRONG SUMMARY) *The strong summary f_s is an RDF summary based on the strong equivalence \equiv_s .*

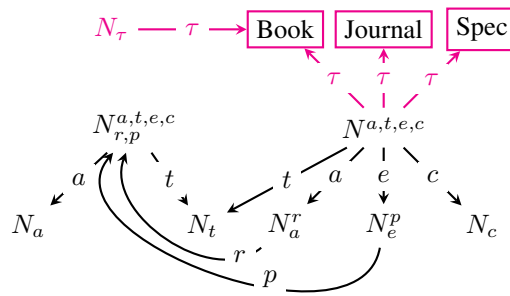


Figure 10: Strong summary of the RDF graph in Figure 3.

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 10. 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 G 's schema.

Compared with the weak summary (Figure 5), the strong summary refines (splits) the weak summary node $N_{r,p}^{a,t,e,c}$ into two nodes, $N_{r,p}^{a,t,e,c}$ and $N^{a,t,e,c}$, because the resources from G having the output clique $\{a, t, e, c\}$ have either an empty target clique, or the target clique $\{r, p\}$. As a consequence, a strong summary may have several edges with the same label: in Figure 10, an a -labeled edge exits $N_{r,p}^{a,t,e,c}$ and another one exits $N^{a,t,e,c}$. In contrast, in a weak summary, only one edge can have a given label (Property 7).

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 .

The following properties can be easily established:

Proposition 11 (CLIQUE NODE PROPERTIES) *Let $N(TC_1, SC_1), N(TC_2, SC_2)$ be distinct nodes in a clique summary S_G .*

1. If $TC_1 = TC_2$, then $SC_1 \neq SC_2$ (and the symmetric, exchanging source and target cliques)
2. If $TC_1 \neq TC_2$, then $TC_1 \cap TC_2 = \emptyset$
3. The data properties of which $N(TC_1, SC_1)$ is a target in S_G form a subset of TC_1 (and similarly, the data properties of which $N(TC_1, SC_1)$ is a source form a subset of SC_1).
4. $N(TC_1, SC_1)$ is the source of two distinct data properties p_1, p_2 in S_G iff there exist in G two resources r_1, r_2 such that the following conditions hold:
 - TC_1 is the target clique of r_1 and also of r_2 ;
 - SC_1 is the source clique of r_1 and also of r_2 ;
 - r_1 is the source of p_1 while r_2 is the source of p_2 .
5. Two properties p_1, p_2 are source-related in S_G iff they are source-related in G .
6. The set of source cliques of S_G is exactly the set of source (resp. target) cliques of G .

Proof 13 *Items 1 to 4 are shown in a straightforward manner.*

Proof of item 5:

“If”:

This is easy to see based on the observation that all the properties of which resource $r \in G$ is a source (target) carry over to its representative $(TC(r), SC(r))$ in S_G . Thus, any resource having two properties in G is represented by a node having the same two properties in S_G , thus source- and target-relatedness is preserved by clique summarization.

“Only If”: *The proof is by induction over k , the distance between the p_1 and p_2 , source-related in S_G .*

First consider $k = 0$: this corresponds to the case when p_1 and p_2 have a common source in S_G . By item 3., resources r_1 and r_2 are represented by the node $N(TC_1, SC_1)$, which is a common source for p_1 and p_2 .

Now consider p_1, p_2 are source-related at distance $k + 1$ in S_G . Then there exists a property p_3 and a node N_3 in S_G such that p_1 and p_3 are source-related at distance k , while N_3 has p_3 and p_2 . By the induction hypothesis, then, p_1 and p_2 are source-related in G ; further, for N_3 to exist in S_G , it must be the case that a G node has p_3 and p_2 , making p_1 and p_2 source-related in G .

The proof of item 6 follows easily from item 5.

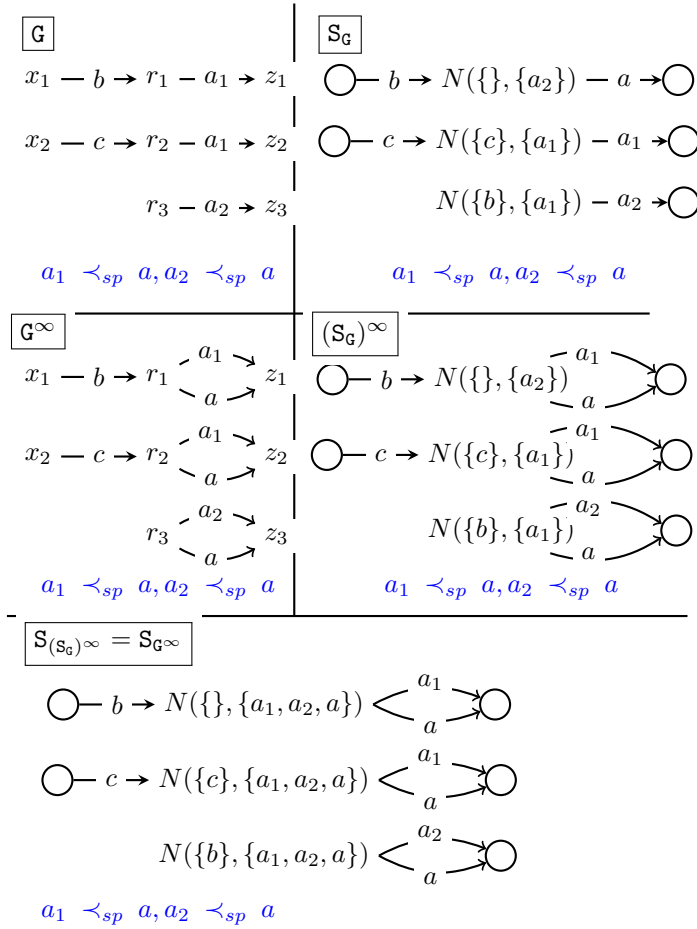


Figure 11: Illustration of the strong completeness statement (some summary nodes are shown as unlabeled circles).

Proposition 12 (STRONG COMPLETENESS) *For a given graph G , $S_{G^\infty} = S_{(S_G)^\infty}$ holds.*

Proof 14 (Data triples)

To prove that the strong summaries of G^∞ and $(S_G)^\infty$ have the same data components, we start by showing that the nodes in the data components of these cliques are obtained by identical calls to N . This reduces to showing the following two statements:

1. two G properties appear in the target set of a given node in S_{G^∞} iff they appear in the target set of a given node in $S_{(S_G)^\infty}$ (and similarly for the source sets);
2. a target and a source clique of G are in the same data node in S_{G^∞} iff they are in the same data node in $S_{(S_G)^\infty}$.

Together, 1. and 2. above ensure that the data properties of G are “grouped together identically” in source, respectively target sets, when building the nodes of the two strong summaries. Both 1. and 2. are implied (are particular cases of) the corresponding statements in the proof of Proposition 8. This is because any two properties p_1, p_2 target-related in $(S_G)^\infty$ are also target-related in $(W_G)^\infty$.

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 relatedness. Relying on the notions of untyped resources UN_G , untyped data graph UD_G and the class set representation function $C(X)$ introduced there, we similarly define:

Definition 16 (U-STRONG EQUIVALENCE AND SUMMARY) *Two nodes in a graph G are untyped-strong equivalent, denoted $n_1 \equiv_{\text{US}} n_2$, 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 US_G , is its summary through \equiv_{US} .

Untyped strong equivalence restricts strong equivalence to untyped resources only. The untyped-strong summary US_G summarizes UD_G in strong fashion, and leaves all typed resources untouched. We can now define:

Definition 17 (TYPED STRONG SUMMARY) *The typed strong summary TS_G of an RDF graph G is the untyped-strong summary of the type-based summary of G , namely US_{T_G} .*

In our example, it turns out that the type-strong summary of the RDF graph in Figure 3 coincides with the type-weak one shown in Figure 8. As can be easily seen from their definition, 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 13 (TYPED STRONG FIXPOINT) *The typed strong summary (Definition 17) has the fixpoint property.*

The proof follows again by from the fixpoint property of the type-based summary, and that of the untyped strong one.

Proposition 14 (TYPED STRONG NON-COMPLETENESS)

There exists an RDF graph G such that $\text{TS}_{G^\infty} \neq \text{TS}_{(\text{TS}_G)^\infty}$ holds.

More generally, in the presence of domain (\leftarrow_d) or range (\leftarrow_r) RDF schema rules, one cannot compute TS_{G^∞} from G because the typed weak summary represents typed resources differently from the untyped ones. The \leftarrow_d and \leftarrow_r schema rules may turn untyped resources into typed ones, thus leading to divergent representations of G 's data nodes in TS_G and respectively TS_{G^∞} .

6 Implementation & experiments

We implemented our summarization tool in Java 1.8 (approximately 15.000 lines of code). It issues SQL queries to the underlying RDF data store implemented in PostgreSQL, and builds the respective summaries based on the results of these queries. We chose a simple, generic, yet robust storage of RDF graphs within the PostgreSQL v9.3.2 relational database server.

The triples are loaded from a file to the triples table in Postgres through the Postgres `COPY` command (currently, only files in n-triples format are supported). The triples table comprises three columns, for subject, property and object. We encode the triples table and subsequently work only with the *integer representation of the input RDF graph*. Operating on integers instead of strings provides for savings both in processing time and memory. For each resource from G , the dictionary table in Postgres stores its unique integer value.

The encoded triples table is split into a table of data triples and a table of type triples, where each row is assigned its sequence number.

We rely on the Jena RDF API [21] to parse the RDFS triples. Typically there are much less schema triples compared to instance triples, so we load the schema directly to memory, after which we encode and store it to a separate Postgres table.

6.1 Data structures

Summary. Similar to the input graph, we work with *an integer representation of the summary resources*, i.e., each URI in the summary is represented by an integer during the summary construction; recall that the summary does not contain literals, and represents literals from the input graph by summary nodes given fresh URIs. We decode the summary to retrieve the property URIs after it is fully built.

Currently, our summary graphs are built in memory, based on the Trove library², which improves on the performance of java.util Object-based collections w.r.t. time and memory by using primitives instead of wrapper classes (`int` in place of `Integer` etc.). However, our algorithms can be adapted to manipulate more scalable (disk-based and/or distributed) data structures, and we are currently working in that direction.

Data nodes. The representation of each data node (whether typed or untyped) from G with exactly one data node in H_G is stored in a `TIntIntMap` map, denoted `rd`. Further, a `TIntObjectMap<TIntSet>` multi-map contains for each data node in H_G a set of represented data nodes from G ; we denote it `dr`.

Data node class set. A typed data node has exactly one (explicit) class set, which we store in a `TIntObjectMap<TIntSet>`, denoted `dcIs`. Observe that this in fact models the summary type edges; each key in the map represents a subject of a type triple and the object is the class from its respective class set.

Data node properties. For mapping data properties to their sources (targets, respectively) we use a `TIntIntMap`, denoted `dpSrc` (`dpTarg`, respectively). For each source (target) data node, a set of its adjacent data properties is stored in a `TIntObjectMap<TIntSet>`, denoted `srcDps` (`targDps`). These structures will be used in weak and typed weak summaries, to facilitate merging of data nodes. Observe that in TW_G only untyped data nodes may be merged, so the typed data nodes of TW_G will not be stored in these structures. On the other hand, during the summarization of data triples in W_G all of them are untyped at that point since the typing is yet unknown, so all the data nodes attached to some data properties in W_G will be stored in the described structures.³

Data edges/triples. We denote with `dtp` the mapping of data properties to data triple(s) they appear in.

Weak: A `TIntObjectMap` maps each data property from W_G to a summary data triple represented by a `DataTriple` Java object. Observe that a distinct data property may appear in only one data triple in W_G .

Typed weak & strong: A Trove map `TIntObjectMap<Collection<DataTriple>>` stores, for each data property from TW_G and S_G , the collection of data triples it appears in. Thus, a data property may appear in multiple data triples of TW_G and S_G , depending on the type of the subject and object of the considered data triple.

Representing class sets. Any class set from G is represented by exactly one data node in TW_G and this information is stored in a `TObjectIntMap<TIntSet>`, denoted `clsd`.

Cliques. All source/target property cliques of a strong summary are maintained in two `List<TIntSet>` lists, denoted `sc` and `tc`, respectively.

Clique IDs: As an integer ID of each clique, we take its index in the respective lists (`sc` and `tc`), incremented by one.

Mapping G data nodes to cliques. Two `TIntIntMap` maps store the mapping of each G subject to the ID of its source clique, and of each G object to the ID of its target clique. These two maps are denoted `sToSc` and `oToTc`, respectively.

²<http://trove.starlight-systems.com/>

³In both cases, as the choice of the structures shows, there can be only one untyped source and one untyped target data node per distinct data property, while any untyped data node may be attached to multiple data properties.

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

```

1:  $data \leftarrow \text{EVAL}(\text{SELECT } s, p, o \text{ FROM } D_G)$ 
2: for each  $s, p, o$  in  $data$  do
3:    $src \leftarrow \text{GETSOURCE}(s, p, W_G)$ 
4:    $targ \leftarrow \text{GETTARGET}(o, p, W_G)$ 
5:   // GETTARGET may have modified  $src$  and vice-versa
6:    $src \leftarrow \text{GETSOURCE}(s, p, W_G)$ 
7:    $targ \leftarrow \text{GETTARGET}(o, p, W_G)$ 
8:   if  $\text{!EXISTS DATATRIPLE}(W_G, src, p, targ)$  then
9:      $\text{CREATEDATATRIPLE}(W_G, src, p, targ)$ 
10:  end if
11: end for
12: procedure  $\text{CREATEDATATRIPLE}(W_G, src, p, targ)$ 
13:    $ntp.put(p, src, p, targ)$ 
14:    $dpSrc.put(p, src)$ ,  $srcDps.put(src, p)$ 
15:    $dpTarg.put(p, targ)$ ,  $targDps.put(targ, p)$ 
16: end procedure

```

Mapping cliques to S_G 's data nodes. Finally, each clique is represented by a data node in S_G and this mapping of a clique ID to its representative data node is stored in two $\text{TIntObjectMap}\langle\text{TIntSet}\rangle$ maps, denoted $scToSrc$ and $tcToTarg$.

For the typed strong summary cliques are computed only for untyped data nodes.

Figure 12 gives an overview of the data structures used in each kind of summary. There is significant sharing and reuse, given the common features shared across the algorithms.

	Weak	Typed Weak	Strong	Typed Strong
rd, dr	✓	✓	✓	✓
$dcls$	✓	✓	✓	✓
$clsd$		✓		✓
$dpSrc, dpTarg$	✓	✓		
$srcDps, targDps$	✓	✓		
ntp	✓	✓	✓	✓
sc, tc			✓	✓
$sToSc, oToTc$			✓	✓
$scToSrc, tcToTarg$			✓	✓

Figure 12: Overview of the data structures for different summaries

Dictionary. The Postgres dictionary table comprises pairs of integer keys and string values, for all encoded resources from G and H_G . The dictionary lookup is necessary on two occasions: (i) when splitting the encoded triples table into encoded data and type tables, we need the integer keys for RDF type and RDFS properties; (ii) to decode the summary after summarization we need the string values for all class and property keys, and possibly for decoding the resources represented by each summary data node. For (i) we load only the necessary entries from the dictionary table, while in (ii) we perform the summary decoding completely in Postgres (via joins with the dictionary table).

6.2 Algorithms

In weak and strong summaries we summarize the data triples and then the type triples. In the typed weak summary though, we first summarize the type triples and then the data triples.

Data triples are read one by one, their subject and object are represented with source and target data nodes, possibly unifying the source and target nodes based on the information newly found, so as to ensure the existence of an homomorphism from G into the summaries (and in particular, that a given G node is consistently mapped to the same summary node, even though it may participate to several triples scattered over the graph).

Summarizing data triples in W_G . 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 W_G .

Algorithm 2 shows how we map the subjects of data triples in G to data nodes in W_G . 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 6-10). 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 (lines 12-13). When they differ, we have to merge them (line 15).

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.

Summarizing type triples in W_G . Algorithm 3 shows the implementation of weak summarization of type triples. First, we retrieve subjects and classes of all type triples (line 1) and we try to represent each resulting pair (lines 4-9) as follows. We look up the data node src representing s (line 14). 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 (line 14) and we add the class to the class set of src (line 18).

Alternatively, if src is empty in line 15, we know that s is *typed-only*, so we add it to the list of typed-only resources, and its class to the list of typed-only classes and leave it unrepresented for the time being (lines 6-8).

When all subjects attached to some data properties in G have been represented, if there are any typed-only resources, we represent them as well (line 11). The procedure `REPRESENTTYPEONLY` (line 21) creates a single data node d representing all resources from the list of typed-only resources and having *all* the types of typed-only resources.

Algorithm 2 Representing the subject s of a data property p in G with a data node in W_G

Input: s, p, W_G

Output: Data triples represented in W_G

Variables:

src_u - data node representing an untyped source of p

src_s - data node representing a (possibly typed) resource s

```

1: procedure GETSOURCE( $s, p, W_G$ )
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( $W_G, s$ )
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  $src_s = src_u$  then
13:      return  $src_s$ 
14:    else
15:      return MERGEDATANODES( $W_G, src_s, src_u$ )
16:    end if
17:  end if
18: end procedure
19: procedure CREATEDATANODE( $H_G, r$ )
20:   $d \leftarrow NEWINTEGER()$ 
21:   $rd.put(r, d), dr.put(d, r)$ 
22:  return  $d$ 
23: end procedure

```

6.2.1 Refined summary

Summarizing data triples in TW_G . The basic procedure for summarizing data triples in TW_G is very similar to the one of weak presented in Algorithm 1, with the difference in lines 14-15: the entries are stored to the maps only if src and $targ$ respectively, are untyped⁴. Further, the mapping of data nodes in TW_G (Algorithm 4) 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.

⁴Since in TW_G 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 W_G **Input:** Type triples table T_G , the summary W_G **Output:** Type triples represented in W_G

```

1:  $typ \leftarrow \text{EVAL}(\text{SELECT } s, c \text{ FROM } T_G)$ 
2:  $toRes \leftarrow []$ ;  $toCls \leftarrow []$ 
3: for each  $(s, c)$  in  $typ$  do
4:    $repr \leftarrow \text{REPRESENTTYPETRIPLE}(W_G, s, c)$ 
5:   if  $repr = false$  then
6:      $toRes.add(s)$ ,  $toCls.add(c)$ 
7:   end if
8: end for
9: if  $toRes.size > 0$  then
10:   $\text{REPRESENTTYPEDONLY}(W_G, toRes, toCls)$ 
11: end if
12: procedure  $\text{REPRESENTTYPETRIPLE}(W_G, s, c)$ 
13:   $d \leftarrow rd.get(s)$ 
14:  if  $d = \perp$  then
15:    return false
16:  end if
17:   $cls_d \leftarrow dcls.get(d)$ ,  $cls_d.add(d, c)$ 
18:  return true
19: end procedure
20: procedure  $\text{REPRESENTTYPEDONLY}(W_G, toRes, toCls)$ 
21:   $d \leftarrow \text{NEWINTEGER}()$ 
22:  for each  $r$  in  $toRes$  do
23:     $rd.put(r, d)$ ,  $dr.put(d, r)$ 
24:  end for
25:   $cls_d \leftarrow dcls.get(d)$ 
26:  for each  $c$  in  $toCls$  do
27:     $cls_d.add(c)$ 
28:  end for
29: end procedure

```

Summarizing type triples in TW_G . 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 (lines 14-15). 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 H_G 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 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.

Algorithm 4 Representing the subject of a data property in G with a data node in TW_G

Input: The subject s of the data triple, the property p of s , the summary TW_G

Output: Data triples represented in TW_G

```

1: procedure GETSOURCE( $s, p, TW_G$ )
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( $TW_G, s$ )
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( $TW_G, src_s, src_u$ )
16:    end if
17:  end if
18: end procedure

```

Algorithm 5 Summarizing type triples in TW_G

Input: Type triples table T_G , the summary TW_G

Output: Type triples represented in TW_G

```

1:  $typ \leftarrow EVAL(SELECT s, c FROM T_G ORDER BY s)$ 
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:      $d \leftarrow CREATEDATANODE(TW_G, cls_s, s)$ 
9:      $clsd.put(cls_s, d)$ 
10:  end if
11: end for

12: procedure CREATEDATANODE( $TW_G, r, cls_r$ )
13:   $d \leftarrow NEWINTEGER()$ 
14:   $rd.put(r, d), dr.put(d, r), dcls.put(d, cls_r)$ 
15:  return  $d$ 
16: end procedure

```

6.2.2 Strong summary

Summarizing data triples in S_G . 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 distinct data node d in G , a representative data node d_{repr} is assigned in S_G . 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 S_G . The method for creating the data triple is the same as for the weak summary in Algorithm 1.

Algorithm 6 Summarizing data triples in S_G

Input: Data triples table D_G , the summary S_G

Output: Data triples represented in S_G

Variables:

$scId, tcId$ - source (target) clique ID

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

```

1: BUILDSOURCECLIQUES( $D_G, S_G$ )
2: BUILDTARGETCLIQUES( $D_G, S_G$ )
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:    $d_{repr} \leftarrow FIRSTCOMMONELEMENT(D_{SC}, D_{TC})$ 
8:   if  $d_{repr} \neq \perp$  then
9:      $rd.put(r, d_{repr}), dr.put(d_{repr}, r)$ 
10:  else
11:     $d_{repr} \leftarrow CREATEDATANODE(S_G, d)$ 
12:     $scToSrc.put(scId, d_{repr})$ 
13:     $tcToTarg.put(tcId, d_{repr})$ 
14:  end if
15: end for
16:  $data \leftarrow EVAL(SELECT \ s, p, o \ FROM \ D_G)$ 
17: for each  $s \ p \ o$  in  $data$  do
18:    $src \leftarrow rd.get(s), targ \leftarrow rd.get(o)$ 
19:   if  $\exists$   $!EXISTSDATATRIPLE(S_G, src, p, targ)$  then
20:      $CREATEDATATRIPLE(S_G, src, p, targ)$ 
21:   end if
22: end for

```

The procedure BUILDSOURCECLIQUES(D_G, S_G) invokes the procedure COMPUTESOURCECLIQUE(S_G, s, P_s) shown in Algorithm 7 for each distinct subject in D_G and the set of data properties of which it is the source.

First, in Algorithm 7 (lines 3-9) 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 (lines 10-11). 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 S_G (lines 13-15). Finally, $srcClq$ is assigned to s in $sToSc$ (line 17).

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 S_G . The types in S_G are summarized in the same manner as for the W_G , described in Algorithm 3.

Algorithm 7 Computing source cliques in S_G

Input: The summary S_G , 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 S_G populated

```

1: procedure COMPUTESOURCECLIQUE( $S_G, s, P_s$ )
2:    $srcClq \leftarrow \perp$ 
3:   for each  $c$  in  $sc$  do
4:      $common \leftarrow \text{FINDFIRSTCOMMONELEMENT}(c, P_s)$ 
5:     if  $common > 0$  then
6:        $srcClq \leftarrow clique$ 
7:       break;
8:     end if
9:   end for
10:  if  $srcClq \neq \perp$  then
11:     $srcClq.addAll(P_s)$ 
12:  else
13:     $srcClq \leftarrow \text{new TIntHashSet}$ 
14:     $srcClq.addAll(P_s)$ 
15:     $sc.add(srcClq)$ 
16:  end if
17:   $sToSc.put(s, sc.INDEXOF(srcClq)+1)$ 
18: end procedure

```

6.3 Complexity of our summary building algorithms

We report here on the I/O complexity incurred by our summary building algorithms, as this is the most significant component of the building time.

Weak & typed weak summaries. The I/O cost (data read by our algorithm in order to build the summary) is: $|D_G| + |T_G|$.

Strong summary. The I/O cost is distributed as follows:

- finding source cliques: $|D_G|$
- finding target cliques: $|D_G|$
- finding all distinct data nodes in G : $2 \times |D_G|$
- one pass to represent data nodes of each data triple: $|D_G|$
- types: $|T_G|$

Thus, the total I/O cost is: $5 \times |D_G| + |T_G|$.

Typed strong summary. The I/O cost is incurred by:

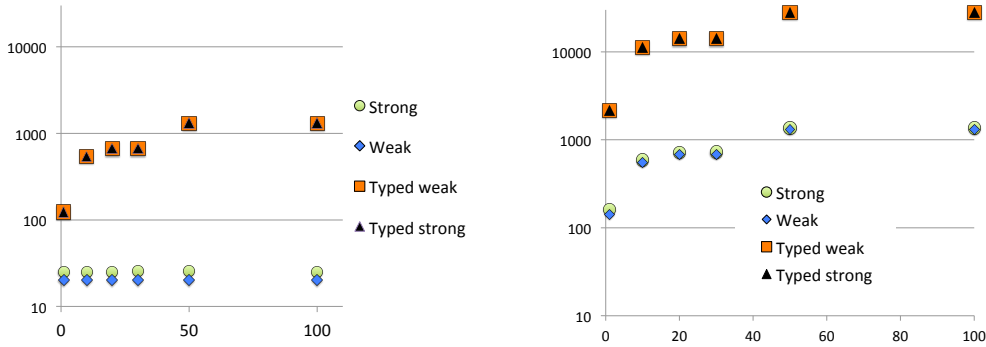


Figure 13: The numbers of data nodes (left) and all nodes (right) in BSBM summaries.

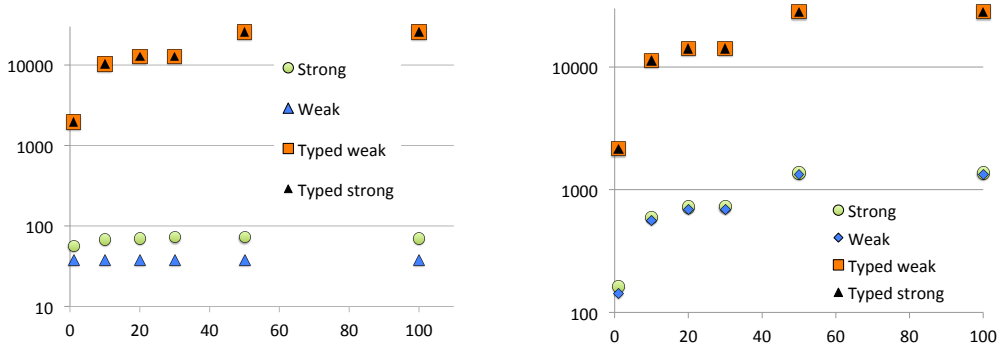


Figure 14: The numbers of data edges (left) and all edges (right) in BSBM summaries.

- reading all the types: $|T_G|$
- computing the source and target cliques: $2 \times |D_G| + 2 \times |T_G|$
- computing the distinct data nodes in G : $2 \times |D_G|$
- making one pass to represent data nodes of each data triple: $|D_G|$

Thus, the I/O cost sums up to $4 \times |T_G| + 4 \times |D_G|$.

6.4 Experimental results

The hardware specification of the machine used for running the summarization tool is the following: Intel Xeon X5647 Processor (2.93GHz, 4 Cores, 8 CPUs, 12M Cache, 5.86 GT/s QPI, 130W TDP, Turbo, HT), DDR3-1066MHz; 16Gb of memory (4x RDIMM 4Go Dual Rank LV) 1066MHz, with modules DIMM 1333MHz; 2 hard disks Hot Plug 600Go SAS 6Gbit/s 15000tr/min 3,5 inches.

The PostgreSQL version 9.3.2 was used and configured as shown in Figure 15. We have run the summarization tool for various JDBC fetch sizes and dataset sizes and have settled on 100,000 as the optimal fetch size for our hardware setting. The experiments were run with 16 GB of RAM assigned to the Java Virtual Machine.

We generated the weak, strong, typed weak and typed strong summary for The Berlin SPARQL Benchmark [2] (BSBM) dataset of various sizes. Figure 13 reports the number of data nodes and the number of all the nodes, respectively, in the four summaries. The horizontal axis is numbered in millions

shared_buffers	4 GB
work_mem	64 MB
maintenance_work_mem	256 MB
checkpoint_segments	256 MB
checkpoint_timeout	1h

Figure 15: The PostgreSQL 9.3.2 configuration

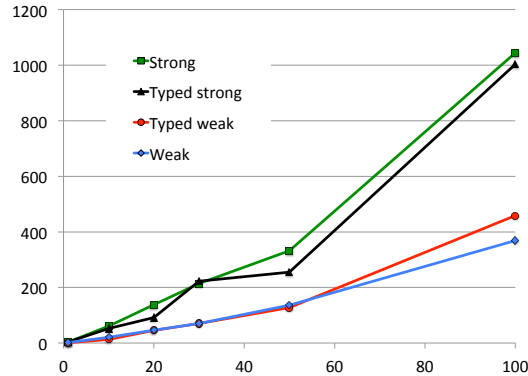


Figure 16: Summarization time (seconds) for various data sizes

of triples in the input, while the vertical axis is in logarithmic scale. The number series are ordered so as to improve readability; observe that S_G and W_G numbers are very close to each other, and so are the TW_G and TS_G numbers (the respective dots sometimes overlap). This shows that in practice there is not a big difference between sharing an input/output clique directly, as nodes are required to do in order to be summarized together in S_G , and through a chain of cliques, as it is the case in W_G . Isolating typed data nodes from the others, as we do in TW_G and TS_G , does have a strong impact, multiplying the number of data nodes by a factor between 5 and 50. Further, for BSBM data, the number of class nodes (the difference between the two numbers recorded in 13) is significantly higher than the number of data nodes (by a factor of 5 up to 60) for the strong and weak summaries, demonstrating the very high summarization power (reduction in data nodes) of these summaries. The number of class nodes goes approximately from 100 to 1300 for all summary sizes.

Figure 14 presents the number of data edges, and the number of total edges in each summary. The figure confirms that the summaries W_G and S_G , which we call *type-first* summaries (to highlight the importance they give to RDF types) behave similarly, while the typed ones are most complex but again similar to each other. The overall number of edges remains very moderate (at most 28210) remains very small compared to the data size (10 to 100 million triples); thus, the summary occupies at most 0.028 of the data size, and in the best case, only 2.8×10^{-4} of the data size.

Figure 16 depicts the time needed to build the summaries using our Postgres-based algorithms. The weak and strong summaries, are built in at most 8 minutes, whereas TS_G takes up to 1000 s (16 minutes) and TW_G up to 32 minutes. We find these times acceptable for a centralized implementation based on a DBMS, especially considering that summarization is an off-line task. The building time increases with the data size, and is higher for the strong and typed strong summaries than for the other two. This is because building strong summaries also requires actually computing the cliques, whereas for the weak ones, this is not needed.

Summary building times and statistics for other datasets. We include below the similar set of metrics for two other data sets we have worked with: John Peel BBC Radio Sessions⁵ of $\approx 270k$ triples and

⁵<http://dbtune.org/bbc/peel/>

INSEE Geo⁶ of $\approx 370k$ triples.

Dataset	Data edges	Type edges	Data nodes	Classes	Total nodes
Peel	24	9	7	9	16
INSEE Geo	21	9	6	9	15

Figure 17: The W_G sizes for various RDF datasets

Dataset	Data edges	Type edges	Data nodes	Classes	Total nodes
Peel	79	20	25	9	34
INSEE Geo	26	10	8	9	17

Figure 18: The S_G sizes for various RDF datasets

Dataset	Data edges	Type edges	Data nodes	Classes	Total nodes
Peel	42	9	17	9	26
INSEE Geo	68	11	15	9	24

Figure 19: The TW_G sizes for various RDF datasets

7 Related Work

OEM and XML summaries. Summaries have long been studied in the context of semistructured data. Dataguides [7] were introduced to summarize semistructured OEM graphs, similar to RDF, but assumed to have a “root” node, from which all others are accessible; this may not hold for RDF. A Dataguide features *exactly once each path in the original graph*, and *each Dataguide path corresponds to one path in the graph*. This allows for *several* distinct dataguides, whose construction from the graph is shown [14] to amount to constructing a deterministic finite automaton out of a non-deterministic one; Dataguides construction has worst-case exponential time complexity, thus is not in general feasible. An algorithm is provided for building *strong* Dataguides, used as a basis for indexing. The 1-index [13] groups together OEM or XML nodes that are reachable by exactly the same set of paths. Later works focused on indexes for supporting XML path queries [4, 9], or path-based XML summarization into graphs [5]. All these works differ from ours, because the input is a tree or DAG and/or because it lacks types and implicit information.

Graph summarization. Graph summarization has been very intensively studied, in particular through mining or clustering; large-scale graph processing is also a hot topic. A large number of works build on the idea of Dataguides for graph data, oftentimes referred to as *structural indexes*, which bear a similarity to graph summaries, both being a reduced version of the input graph and collapsing nodes based on some common attributes.

Our focus is on *RDF graphs* with *implicit* data, for which we devised *query-oriented* summaries, which are RDF graphs themselves and may be computed on a variety of platforms.

Graph *cores* have been studied in [6]. 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 property. Our summaries are not cores of the incoming graph G , since we cannot guarantee a homomorphism from either summary to the graph G . In exchange, both summary versions we consider can be built in polynomial time in the size of G , while computing the core is much harder.

Bisimulation-based approaches group nodes by the similarity of their neighborhood [11, 17]. In [11] bisimulation is used to summarize graphs from the LOD cloud, focusing on the distribution of classes and properties across LOD sources. The resulting *resource-oriented* summary comprises unlabeled edges. [17] utilizes bisimulation to construct a structural index from structure patterns exhibiting certain edge

⁶<http://rdf.insee.fr/geo/index.html>

Dataset	Data edges	Type edges	Data nodes	Classes	Total nodes
Peel	52	9	21	9	30
INSEE Geo	69	11	16	9	25

Figure 20: The TS_G sizes for various RDF datasets

Dataset	Input size	Weak	Strong	Typed Weak	Typed Strong
Peel	271369	0.307	1.137	0.283	0.842
INSEE Geo	368761	0.360	1.447	0.348	1.130

Figure 21: The summarization times in seconds for various RDF datasets

labels that comprise paths of some maximum length. The main problem with bisimulation is that as the size of the neighborhood increases, the size of bisimulation grows exponentially and can be as large as the input graph. Thus, as we aim for both complete and compact summaries, bisimulation is not a good fit.

To overcome the problem with bisimulation, [8] suggests locality-based summaries, generated by a graph partitioning algorithm. Nonetheless, a reduction of the input RDF graph is necessary which is achieved by removing triples having properties with literal values, thus resulting in an incomplete summary. Recall that we wish to preserve queries comprising properties with literal values as well.

A *triple-oriented* structural index for RDF data is built in [16] as a non-RDF graph, where a node comprises a set of triples forming a data partition, while an edge describes the way in which triples from its adjacent nodes join. [18] proposes a tree RDF index, storing regions defined by center vertices, limited to property paths and built assuming that the input RDF graph is saturated.

In [10], RDF classes are inferred based on the common properties of resources. Thus, only common source patterns are analyzed, while the common targets and property paths are not considered. Further, the `rdf:type` properties that a dataset may comprise are simply ignored.

[3] explores alternative RDF summaries w.r.t. graph homomorphism and trades precision for efficiency in computing them.

Summarizing implicit data is not considered in these works.

In [15], the authors consider building compact summaries of RDF databases, and they take into account semantic knowledge so that minimal summaries are extracted. Specifically, a summary should account in a minimal way for information present in an RDF dataset, in the presence of subclass constraints only. The authors do not consider the other classes of rules that we consider in our work (subproperty, domain and range) and do not study the formal summary properties we identify here.

References

- [1] M. Arias, J. D. Fernández, M. A. Martínez-Prieto, and P. de la Fuente. An empirical study of real-world SPARQL queries. *CoRR*, abs/1103.5043, 2011.
- [2] C. Bizer and A. Schultz. The berlin SPARQL benchmark. *Int. J. Semantic Web Inf. Syst.*, 5(2):1–24, 2009.
- [3] S. Campinas, R. Delbru, and G. Tummarello. Efficiency and precision trade-offs in graph summary algorithms. In *IDEAS*, 2013.
- [4] Q. Chen, A. Lim, and K. W. Ong. $D(K)$ -index: An adaptive structural summary for graph-structured data. In *SIGMOD*, 2003.
- [5] M. P. Consens, R. J. Miller, F. Rizzolo, and A. A. Vaisman. Exploring XML web collections with DescribeX. *ACM TWeb*, 4(3), 2010.

-
- [6] C. Godsil and G. Royle. *Algebraic graph theory*. Springer-Verlag, 2001.
 - [7] R. Goldman and J. Widom. Dataguides: Enabling query formulation and optimization in semistructured databases. In *VLDB*, 1997.
 - [8] S. Gurajada, S. Seufert, I. Miliaraki, and M. Theobald. Using graph summarization for join-ahead pruning in a distributed RDF engine. In *SWIM Workshop*, 2014.
 - [9] R. Kaushik, P. Bohannon, J. F. Naughton, and H. F. Korth. Covering indexes for branching path queries. In *SIGMOD*, 2002.
 - [10] K. Kellou-Menouer and Z. Kedad. A clustering based approach for type discovery in RDF data sources. In *Extraction et Gestion des Connaissances*, 2015.
 - [11] S. Khatchadourian and M. P. Consens. ExpLOD: Summary-based exploration of interlinking and RDF usage in the linked open data cloud. In *ESWC*, 2010.
 - [12] D. Lanti, M. Rezk, G. Xiao, and D. Calvanese. The NPD benchmark: Reality check for OBDA systems. In *EDBT*, 2015.
 - [13] T. Milo and D. Suciu. Index structures for path expressions. In *ICDT*, 1999.
 - [14] S. Nestorov, J. D. Ullman, J. L. Wiener, and S. S. Chawathe. Representative objects: Concise representations of semistructured, hierarchical data. In *ICDE*, 1997.
 - [15] M. Palmonari, A. Rula, R. Porrini, A. Maurino, B. Spahiu, and V. Ferme. ABSTAT: linked data summaries with abstraction and statistics. In *ESWC 2015 Satellite Events*, pages 128–132, 2015.
 - [16] F. Picalausa, Y. Luo, G. H. L. Fletcher, J. Hidders, and S. Vansummeren. A structural approach to indexing triples. In *ESWC*, 2012.
 - [17] T. Tran, G. Ladwig, and S. Rudolph. Managing structured and semistructured RDF data using structure indexes. *IEEE Trans. Knowl. Data Eng.*, 25(9):2076–2089, 2013.
 - [18] O. Udea, A. Pugliese, and V. S. Subrahmanian. GRIN: A graph based RDF index. In *AAAI*, 2007.
 - [19] W3C. Resource description framework. <http://www.w3.org/RDF/>.
 - [20] <http://gromgull.net/blog/2010/09/btc2010-basic-stats>, 2010.
 - [21] Jena. <http://jena.sourceforge.net>.



**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