



**HAL**  
open science

## Query-Oriented Summarization of RDF Graphs

Šejla Čebirić, François Goasdoué, Ioana Manolescu

► **To cite this version:**

Šejla Čebirić, François Goasdoué, Ioana Manolescu. Query-Oriented Summarization of RDF Graphs. [Research Report] RR-8920, INRIA Saclay; Université Rennes 1. 2017. hal-01325900v5

**HAL Id: hal-01325900**

**<https://inria.hal.science/hal-01325900v5>**

Submitted on 28 Jun 2017 (v5), 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

ISRN INRIA/RR--8920--FR+ENG

ISSN 0249-6399





## Query-Oriented Summarization of RDF Graphs

Šejla Čebirić, François Goasdoué, Ioana Manolescu

Project-Teams CEDAR

Research Report n° 8920 — version 4 — initial version June 2016 — revised  
version June 2017 — 46 pages

**Abstract:** RDF is the data model of choice for Semantic Web applications. RDF graphs are often large and heterogeneous, thus users may have a hard time getting familiar with the structure and semantics of a graph, or determining whether a graph is useful for a certain application.

We consider answering such questions by inspecting a *graph summary*, a compact structure conveying as much information as possible about the input graph. A summary is *representative* of a graph if it represents both its explicit and implicit triples, the latter resulting from RDF Schema constraints. To ensure representativeness, we define a novel RDF-specific summarization framework based on *RDF node equivalence* and *graph quotients*; our framework can be instantiated with many different RDF node equivalence relations. We show that our summaries are representative, and establish a *sufficient condition* on the RDF equivalence relation to ensure that a graph can be *efficiently summarized*, without materializing its implicit triples. We demonstrate that the state-of-the-art *bisimulation* equivalence relations between graph nodes fit into our framework. Further, we instantiate the framework through *four novel summaries*, based on the new concept of *property cliques*, specifically tailored to cope with highly heterogeneous RDF graphs; we show that they are orders of magnitude more compact than bisimulation summaries. Finally, we show that the bisimulation and two of our clique summaries can be built efficiently so that they represent the explicit and implicit data of the input graph without saturating the graph. The performance benefits of our efficient summarization method is confirmed through a set of experiments.

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

RESEARCH CENTRE  
SACLAY – ÎLE-DE-FRANCE

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

## Résumés orientés requêtes de graphes RDF

### Résumé :

RDF est le modèle de données du W3C, fondé sur les graphes, pour les applications du Web Sémantique. Les graphes RDF sont souvent larges et hétérogènes, ce qui complique la tâche des utilisateurs qui tentent de se familiariser avec leurs structure et sémantique, ou de déterminer si un graphe est utile pour une application donnée.

Nous étudions comment répondre à ces besoins en inspectant un *résumé de graphe*, une structure compacte portant autant d'information que possible à propos du graphe d'entrée. Un résumé est *représentatif* d'un graphe s'il représente à la fois ses triplets explicites et implicites, ces derniers découlant de contraintes de RDF Schema. Pour garantir la représentativité, nous définissons un nouveau cadre de résumé de graphe spécifique à RDF, fondé sur une *relation d'équivalence de noeuds RDF* et le *quotient de graphe* ; notre cadre peut être instancié avec une multitude de relations d'équivalence de noeuds RDF. Nous montrons que nos résumés sont représentatifs et nous établissons une *condition suffisante* sur la relation d'équivalence de noeuds RDF pour garantir qu'un graphe peut être *efficacement* résumé, sans matérialiser ses triplets implicites. Nous démontrons que les relations d'équivalence entre noeuds de l'état de l'art, fondées sur la *bissimulation* de graphes, satisfont notre condition. De plus, nous instancions notre cadre avec quatre nouvelles relations d'équivalence de noeuds RDF, donnant lieu à quatre nouveaux types de résumés, fondées sur les *cliques de propriétés*, spécialement adaptées à la grande hétérogénéité des graphes RDF ; nous montrons que ces types de résumé sont plusieurs ordres de magnitude plus compacts que les résumés à base de bissimulation. Enfin, nous montrons que les résumés à base de bissimulation et deux de nos résumés à base de cliques peuvent être construits efficacement, de sorte qu'ils représentent les triplets explicites et implicites du graphe d'entrée sans avoir à le saturer. Les bénéfices de performance de notre méthode de résumé sont confirmés au travers d'un ensemble d'expériences.

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

# 1 Introduction

The Resource Description Framework (RDF) is the W3C standard for Semantic Web applications. RDF graphs are *varied*, coming from scientific applications, social or online media, government data etc. They are often *large*, and *heterogeneous*, i.e., resources described in an RDF graph may have very different sets of properties. An RDF resource may have one or several *types* (which may or may not be related to each other) or lack types. *RDF Schema* (RDFS) statements are sometimes part of an RDF graph, in which case they lead to **implicit data**.

For illustration, consider the RDF graph  $G_1$  consisting only of the triples shown with solid edges in Figure 1; here and throughout the paper, class nodes, property nodes and RDF Schema triples are shown in blue (we formally revisit them in Section 2). The graph describes resources  $p_1, p_2, c_1, o_1, o_2, o_3$  and relationships between them, while the RDF Schema triples state respectively that someone having the property *heldOffice* is a *Politician*, that *havingSharesIn* something is a form of *ownership*, and that something in which one can have shares, is a *Company*. Implicit triples are shown with dotted edges: thus, because  $p_2$  held offices  $o_2$  and  $o_3$ , the graph implicitly states that  $p_2$  is of type *Politician*. Similar reasoning leads to the other dotted (implicit) triples in the figure;  $G^\infty$  consists of all the (solid or dotted) triples. Importantly, **the semantics of an RDF graph  $G$ , typically denoted  $G^\infty$ , comprises both its explicit and implicit data**. In particular, when querying RDF graphs through SPARQL, the W3C’s standard query language, query answers must *reflect all the explicit and implicit data*.

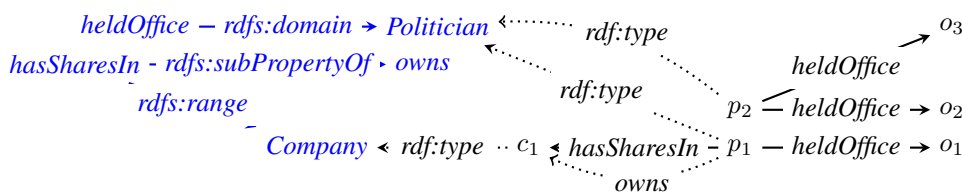


Figure 1: Sample RDF graph  $G_1$  with RDF Schema triples and implicit triples.

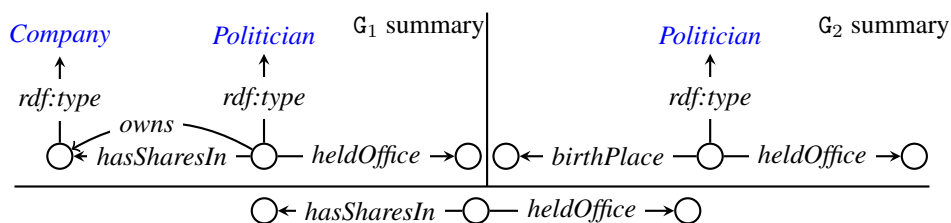


Figure 2: Top left: summary of the implicit and explicit triples of  $G_1$ . Top right: summary of the implicit and explicit triples of another graph  $G_2$ . Bottom: summary of the explicit triples of  $G_1$  only.

The complexity, heterogeneity, and possible semantics associated to an RDF graph may make them difficult to get acquainted with. For instance, consider a journalist interested in *politicians owning companies*. Let  $G_1$  be the graph in Figure 1, and  $G_2$  be another graph specifying for some politicians the offices they held and their birth places. Clearly,  $G_1$  is the best source to use in this case, as  $G_2$  lacks company information. How can we help the journalist figure this out quickly?

One solution is to use graph summaries [12, 11, 9] and in particular those specific to RDF graphs, e.g., [31, 20, 26]. Without loss of generality, a summary is a smaller RDF graph which “represents” the RDF graphs in some sense. Users may consult the summary to get an idea if the original graph is interesting for their needs. For instance, the summaries of  $G_1$  and  $G_2$  could be those shown in Figure 2,

quickly guiding the journalist toward  $G_1$ . In the figure and throughout this paper, circles denote summary nodes, each of which is derived from (represents) a set of nodes from the RDF graph.

If  $G$  is a strict subset of  $G^\infty$ , the summary of  $G$  may be different from that of  $G^\infty$ . However, given that  $G^\infty$  is the semantics of  $G$ , *the summary really representing  $G$  is in fact the one of  $G^\infty$ !* For instance, at the bottom of Figure 2 appears the summary of only the explicit triples of  $G_1$ . This summary does not reflect that in  $G_1$ , *politician  $p_1$  owns part of Company  $c_1$* , as this triple is implicit. In contrast, the summary at top left in Figure 2 fully represents  $G_1$ , that is, it summarizes (represents) its full semantics  $G_1^\infty$ .

Ideally, the summary of  $G_1^\infty$  could be made available by the site hosting  $G_1$ . What if this is not the case? Assuming the owners of  $G_1$  are willing to share it, one could try to saturate it, but RDF/SPARQL endpoints usually do not allow updates, thus one cannot add the implicit triples directly to  $G_1$ . Instead,  $G_1$  could be copied elsewhere, saturated, and  $G_1^\infty$  summarized. This is however quite cumbersome, especially if the saturation is large and/or building it takes a long time. It would be more efficient if one could get the summary of  $G_1^\infty$  directly based on  $G_1$ , without needing to saturate it.

This work is the first study of fully representative RDF graph summarization, taking into account both its explicit and implicit triples. Toward *representative, compact and efficient summarization of RDF graphs*, this work makes the following contributions:

1. We consider RDF graph summarization based on the classical notion of *quotient graphs*: this corresponds to a very wide family of graph summaries, each defined through one notion of *node equivalence*. While numerous useful graph summaries from the literature fit this framework, we show that *a naïve application of quotient summarization to RDF graphs with type and schema information prevents a summary from being representative*.

2. To solve this problem, we define *RDF node equivalence relations* based on the data triples of an RDF graph, and define an *RDF summaries* as a quotient graph based on such a relation. We show that any such summary is *representative*, that is: for any query  $q$  from a rich SPARQL subset  $\mathcal{Q}$ , if  $q$  evaluated on a graph  $G$  has some answers, then  $q$  also has answers on the summary. Conversely, we show that the summary is *accurate*: if  $q$  has answers on the summary, then it also has answers on *some* graph which it summarizes<sup>1</sup>.

3. While the summary of  $G^\infty$  is the one really representing  $G$ , building  $G^\infty$  in order to summarize it may be costly, or impossible for  $G$  users that do not have the right to add triples to  $G$ , e.g., those querying a graph available through a SPARQL endpoint. To efficiently compute representative summaries, we describe a *shortcut* method to compute the summary of  $G^\infty$  *directly from  $G$ , without saturating it*, and we provide a *sufficient condition* on an RDF node equivalence relation, for the corresponding RDF summary to be computable through our shortcut. We show that this condition holds for summaries based on the well-known *bisimulation* [15] notion (revisited in our framework).

4. We define four *RDF summaries* (Section 7 and 8) within our framework, based on the new concept of *property cliques* (Section 6), corresponding to a form of “association” between data properties in RDF graphs; a particularly interesting property of cliques is their tolerance to heterogeneity among graph nodes.

5. From an *efficiency* perspective, we demonstrate that two of these four summaries can be built through the shortcut method which the other two summaries do not admit.

6. We implemented a summarization tool based on bisimulation and clique node equivalence relations, and illustrate the performance benefits of building  $G^\infty$  summaries without saturating  $G$ , for the summaries which admit the shortcut. Our experiments demonstrate the *compactness* of our four clique summaries, *orders of magnitude smaller than bisimulation summaries*.

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

Our summaries are currently being used by the ILDA Inria team (<https://ilda.saclay.inria.fr/>) in an application providing visual entrypoints into the Linked Open Data cloud (<http://wilda.inria.fr/>).

<sup>1</sup>Many RDF graphs may have the same RDF summary.

Assertion	Triple
Class	$s \text{ rdf:type } o$
Property	$s \text{ p } o$

Constraint	Triple
Subclass	$s \prec_{sc} o$
Subproperty	$s \prec_{sp} o$
Domain typing	$p \leftrightarrow_d o$
Range typing	$p \hookrightarrow_r o$

Figure 3: RDF (left) &amp; RDFS (right) statements.

lri.fr:8081/lo datlas/). The application calls our summarization tool to obtain a compact yet representative summary, then applies visualization tools to present the summary to the user in order to guide his search for interesting datasets.

## 2 Preliminaries: RDF Graphs and Queries

An *RDF graph* (or *graph*, in short) is a set of *triples* of the form  $s \text{ p } o$ . A triple states that its *subject*  $s$  has the *property*  $p$ , and the value of that property is the *object*  $o$ . We consider only well-formed triples, as per the RDF specification [32], using uniform resource identifiers (URIs), typed or untyped literals (constants) and blank nodes (unknown URIs or literals).

**Notations.** We use  $s$ ,  $p$ , and  $o$  as placeholders for subjects, properties and objects, respectively. Literals are shown as strings between quotes, e.g., “*string*”. Figure 3 (left) shows how to use triples to describe resources, that is, to express class and property assertions. The RDF standard [32] has a set of *built-in classes and properties*, as part of the `rdf:` and `rdfs:` pre-defined namespaces. We use these namespaces exactly for these classes and properties, e.g., `rdf:type` specifies the class(es) to which a resource belongs. For brevity, we will sometimes use  $\tau$  to denote `rdf:type`.

For example, the RDF graph in Figure 1 describes two politicians identified by the respective URIs  $p_1$  and  $p_2$ , three offices  $o_1$  to  $o_3$ , and a company  $c_1$ .

**RDF Schema (RDFS)** allows specifying *semantic constraints* between the classes and the properties of an RDF graph. Figure 3 (right) shows the four kinds of RDFS constraints, and how to express them through triples. For concision, we denote the properties expressing *subclass*, *subproperty*, *domain* and *range* constraints by the symbols  $\prec_{sc}$ ,  $\prec_{sp}$ ,  $\leftrightarrow_d$  and  $\hookrightarrow_r$ , respectively.

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

For instance, the graph in Figure 1 has the constraints: *heldOffice*  $\leftrightarrow_d$  *Politician*, *hasSharesIn*  $\prec_{sp}$  *owns* and *hasSharesIn*  $\hookrightarrow_r$  *Company*. These constraints lead to the following implicit triples being part of the RDF graph:  $p_1 \tau$  *Politician*,  $p_2 \tau$  *Politician*,  $c_1 \tau$  *Company* and  $p_1$  *owns*  $c_1$ . These triples are shown with dotted lines in the figure.

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



**The semantics of an RDF graph is its saturation.** RDF entailment is part of the RDF standard itself; in particular, *the answers to a query posed on  $G$  must take into account all triples in  $G^\infty$*  [32].

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

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

Further, we call *data property* any property  $p$  occurring in  $D_G$ , and *data triple* any triple in  $D_G$ . For instance, in the RDF graph of Figure 1, data properties are *heldOffice* and *hasSharesIn*, used in the graph’s four data triples.

**2. The graph-based representation of an RDF graph.** As per the RDF specification [32], *the set of nodes* of an RDF graph is the set of subjects and objects of triples in the graph, while its *edges* correspond to its triples. We define three categories of RDF graph nodes: (i) a *class node* is any node whose URI appears as subject or object of a  $\prec_{sc}$  triple, or object of a  $\leftrightarrow_d$  or  $\leftrightarrow_r$  or  $\tau$  triple; (ii) a *property node* is any node whose URI appears as subject/object of a  $\prec_{sp}$  triple, or subject of a  $\leftrightarrow_d$  or  $\leftrightarrow_r$  triple, or property of a triple<sup>2</sup> (iii) a *data node* as any node that is neither a class nor a property node. Note that the sets of class nodes and of property nodes may intersect (indeed, nothing in the RDF specification forbids it), while class nodes are disjoint from data nodes, and property nodes are disjoint from data nodes, too.

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

**Queries.** We consider the SPARQL dialect consisting of *basic graph pattern* (BGP) queries, a.k.a. conjunctive queries, widely considered in research but also in real-world applications. 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  $\{t_1, \dots, t_\alpha\}$  is a BGP; the query head variables  $\bar{x}$  are called *distinguished variables*, and are a subset of those in  $t_1, \dots, t_\alpha$ . For boolean queries,  $\bar{x}$  is empty. The head of  $q$  is  $q(\bar{x})$ , its body is  $t_1, \dots, t_\alpha$ ;  $x, y, z$ , etc. denote variables.

**Query answering.** The *evaluation* of a query  $q$  against  $G$  has access only to the explicit triples of  $G$ , thus it may fail to return the complete answer; the latter is obtained by evaluating  $q$  against  $G^\infty$ . For instance, the query  $q(x_2) :- x_1 \tau \textit{Politician}, x_1 \textit{hasName} x_2, x_1 \textit{owns} c_1$  asks for name of the politician(s) owning company  $c_1$ . Assuming the graph in Figure 1 has an extra triple stating that the name of  $p_1$  is “Cahuzac”, the answer to this query against the explicit and implicit triples of the graph is:  $q(G^\infty) = \{\langle \textit{Cahuzac} \rangle\}$ . Note that evaluating  $q$  only against  $G$  leads to the empty answer.

### 3 Summarization principles

Without loss of generality, we consider that *the summary of an RDF graph  $G$  is an RDF graph itself*. This allows in particular to use existing RDF tools to query or visualize the summary. Further, summarization

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

should satisfy the following conditions:

- *Schema independence*: It must be possible to summarize  $G$  whether or not it has RDF Schema information.
- *Fully automatic*: Summarization should not require user input: the summary should help users learn about the data, not expect them to already know it<sup>3</sup>.
- *Equal treatment of implicit and explicit data*: From the summary of  $G$ , it must be possible to obtain as much information about the explicit triples of  $G$  as about the implicit ones, given that the semantic of  $G$  considers them equally.

Another desirable property is *accuracy*, i.e., the summary should avoid, to the extent possible, reflecting data that does not exist in  $G^\infty$ .

**Criteria for representativeness and accuracy.** We formalize these criteria through the prism of query answering, since RDF graphs are typically accessed and exploited by means of queries. For instance, in our motivating example (Section 1), to find interesting graphs in a large graph collection, one can answer the query “politicians owning companies” over all the graph summaries, and focus on those leading to non-empty answers, e.g.,  $G_1$  in Section 1. Given an *RDF query language (dialect)*  $\mathcal{Q}$ , we define:

**Definition 1.** (REPRESENTATIVENESS) *Let  $G$  be any RDF graph. Its summary  $\text{Sum}$  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  $q(\text{Sum}^\infty) \neq \emptyset$ .*

Informally, queries having answers on  $G$  should also have answers on the summary. This is desirable in order for the summary to help users formulate queries: the summary should reflect all graph patterns that occur in the data.

Formalizing accuracy is a bit more involved. Summarization inevitably leads to some loss of information, which from our query-oriented perspective can be thought of as “false positives” (queries having answers on the summary of  $G$ , but not on  $G$ ) and/or “false negatives” (queries having answers on  $G$ , but not on its summary). As our discussion has highlighted, representativeness is an important asset, thus we rule out false negatives. Therefore, we allow summarization to introduce some false positives, and take this into account when defining accuracy. Note that several RDF graphs may have the same summary. We term *inverse set* of a summary  $\text{Sum}$ , the set of all RDF graphs whose summary is  $\text{Sum}$ ; the inverse set is purely conceptual, i.e., we never actually compute it. Based on this, we define accuracy as:

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

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

Which query dialect should we define representativeness and accuracy for? In order for summaries to be compact, they should not include subject and object values (whether they are URIs, literals or blank nodes); this follows the observation that there are many orders of magnitude more distinct subject and objects, than there are distinct properties in an RDF graph [34]. However, we chose to preserve types (classes) and data properties, as these are essential features of an RDF graph. These choices *rule out representativeness for queries with URIs or literals in subjects or objects positions*. Instead, we identify two large and useful classes of BGP queries:

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

<sup>3</sup>For a different summary usage scenario, it is easy to extend our approach to summarize only triples referring e.g., to some classes and/or properties; we do not pursue this further here.

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

Both languages forbid URIs or literals in subject and object positions, and require that if type triples appear in the query, then the type is known. They differ in that RBGPs require URIs in the property positions, whereas RBGP\*s also allow variables there. Clearly, RBGPs are a restriction of RBGP\*s.

A sample RBGP is  $q(x_2) :- x_1 \tau \text{ Politician}, x_1 \text{ hasName } x_2, x_1 \text{ owns } c_1$  while a similar RBGP\* is  $q^*(x_2) :- x_1 \tau \text{ Politician}, x_1 \text{ hasName } x_2, x_1 y c_1$ .

We define *RBGP\* representativeness* and *RBGP\* accuracy* by instantiating  $\mathcal{Q}$  in Definition 1 and Definition 2, respectively, to RBGP\* queries (Definition 3).

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

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

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

## 4 RDF summarization

As has often been attempted for non-RDF and RDF graphs, we define summaries through the classical notion of *quotient graph* from graph theory, parameterized by a *node equivalence relation*. Section 4.1 recalls these notions and shows that applied as such, they prevent RDF summary representativeness. Section 4.2 presents our approach toward defining RDF summaries that are rich with information and satisfy requirements identified in Section 3, in particular representativeness.

### 4.1 Quotient summaries and their limitations for RDF graphs

Recall that an *equivalence relation*  $\sim$  over a set  $S$  is a binary relation (subset of  $S \times S$ ) which is reflexive, transitive, and symmetric. Further, we recall:

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

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

Many prior graph summarization and indexing works have relied on graph quotients, e.g., [24, 9, 27] (see also Section 12). In particular, these works used the classical notion of *bisimulation* as a node equivalence relation; we recall it below:

**Definition 5.** (FW BISIMILARITY) *Forward (FW) bisimilarity is a node equivalence relation such that two nodes  $n_1, n_2 \in G$  are FW-bisimilar, denoted  $n_1 \sim_{fw} n_2$ , iff (i) for every  $G$  edge of the form  $n_1 \xrightarrow{a} m_1$  there exists a  $G$  edge of the form  $n_2 \xrightarrow{a} m_2$  such that  $m_1 \sim_{fw} m_2$ , and conversely (ii) for every  $G$  edge of the form  $n_2 \xrightarrow{a} m_2$  there exists a  $G$  edge of the form  $n_1 \xrightarrow{a} m_1$  such that  $m_1 \sim_{fw} m_2$ .*

Intuitively, FW bisimilarity requires nodes to have the same sets of outgoing properties, and that these properties lead to respectively bisimilar nodes. The dual relation of forward bisimilarity is *backward (BW) bisimilarity*, noted  $\sim_{bw}$ , which compares nodes based on the incoming edges *only*. Nodes can also be compared based on their incoming *and* outgoing edges using *forward-backward (FB) bisimilarity*, noted  $\sim_{fb}$ . The definitions of these node equivalence relations easily derive from Definition 5; clearly,  $\sim_{fb}$  implies  $\sim_{fw}$  and  $\sim_{bw}$  but the opposites do not hold. For instance, in the graph  $G$  shown Figure 4, it is easy to check that  $u_1 \sim_{fb} u_3, u_2 \sim_{fb} u_4, p_1 \sim_{fb} p_3$  and  $p_2 \sim_{fb} p_4$ .

Bisimilarity essentially requires that the graph structure leading to/starting from two bisimilar nodes  $n_1, n_2$  be the same. A more relaxed condition of *bounded- or  $k$ -bisimilarity* requires that these be equal only up to length  $k$ . For instance, if  $\sim_{fb}^k$  is the bounded version of  $\sim_{fb}$ , and  $n_1 \sim_{fb}^k n_2$ , the graph structures at a distance greater than  $k$  from  $n_1$  and  $n_2$  may differ. Clearly, as bounded bisimilarity is easier to satisfy, a summary based on it will have fewer nodes than one based on pure bisimilarity. Our analysis in this paper applies to bounded and unbounded bisimilarity alike; for space reasons, we will refer only to the latter in the sequel.

Importantly, *straightforward application of quotient summarization on RDF graphs prevents the summary from being representative*. Figure 4 illustrates this on a FB bisimilarity quotient; similar examples can be exhibited for any quotient based on an equivalence relation which ignores the special role types and schema play in RDF graphs. The quotient shown at right in Figure 4 is not a representative summary: for instance, although the query  $q(x) :- x \tau A$  has a non-empty answer on the graph  $G$ , it is empty on the quotient, and the same holds for the query  $q(x) :- x p_2 y$ . This loss of representativity stems from three issues: (•) *loss of class and (some) property names*: the nodes corresponding to classes  $A$  and  $B$  are equivalent and thus fused, since they only have incoming type edges and no outgoing edges; similarly, the property nodes  $p_2$  and  $p_4$  are fused thus they are absent from the summary; (◊) *loss of RDF Schema triples*: the subproperty ( $\prec_{sp}$ ) triples are summarized in a single one which states that “some properties are subproperties of others”; (○) *loss of implicit triple information*: the  $G$  triples shown in dashed edges are not reflected in the quotient, nor can they be inferred from it, given the loss of the schema triples.

## 4.2 Our RDF summarization framework

The first step we take toward RBGP\* representative summaries is to define a generic notion of node equivalence specifically designed for RDF graphs:

**Definition 6.** (RDF NODE EQUIVALENCE) *Let  $\equiv$  be a binary relation between the nodes of an RDF graph. We say  $\equiv$  is an RDF node equivalence relation (or RDF equivalence, in short) iff (i)  $\equiv$  is an equivalence relation (ii) any class node is  $\equiv$  only to itself, and (iii) any property node is  $\equiv$  only to itself.*

With this generic equivalence relation in place, we define our summary of RDF graphs as follows:

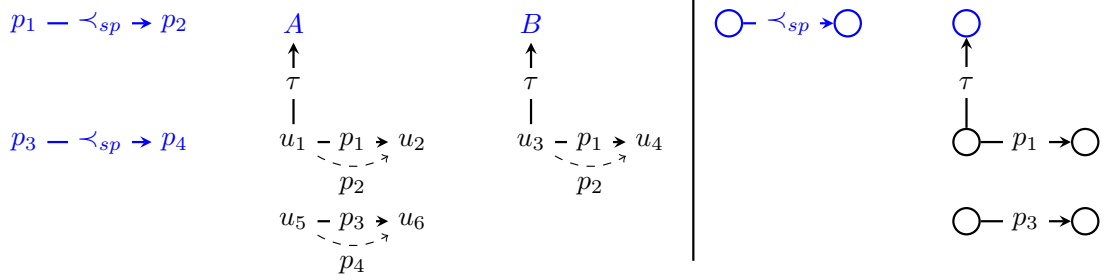


Figure 4: Sample RDF graph  $G$  (left) and its quotient graph based on  $\sim_{fb}$  (right).

**Definition 7.** (RDF SUMMARY) *The summary of an RDF graph  $G$  by an RDF node equivalence relation  $\equiv$ , denoted  $G_{/\equiv}$ , is the RDF graph defined as the quotient graph of  $G$  by  $\equiv$ , in which every class (resp. property) node is named after its  $G$  URI. Further, its data nodes are fresh URIs.*

Importantly, the RDF summaries enjoy the following property:

**Property 1** (Class, property and schema preservation). *An RDF graph  $G$  and an RDF summary  $G_{/\equiv}$  of it have the same sets of classes names, of property names, and of schema triples.*

Indeed, Definitions 6 and 7 ensure that class and property nodes are preserved through summarization, as well as their URI labels, since they cannot be equivalent (hence fused) with other nodes. Further, because our summarization approach relies on graph quotients, all property names labelling edges in a given graph also label edges in its summary. This obviously implies that schema triples are preserved, as they only involve class or property nodes.

It is worth noting that the above property de facto prevents in RDF summaries the issues (●) and (◇) pointed out in the preceding section. However, a solution to the third issue (○) requires the following result:

**Proposition 1.** (SUMMARY REPRESENTATIVENESS) *An RDF summary is RBGP\*-representative.*

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

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

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

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

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

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

- If  $\phi(t')$  is a *schema* triple, then by definition of RDF entailment  $\phi(t') \in (S_G)^\infty$ , thus it is also in  $(G_{/\equiv})^\infty$  since  $S_G = S_{G_{/\equiv}} \subseteq G_{/\equiv}$  by definition of a summary.

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

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

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

- If both  $\phi(t_1)$  and  $\phi(t_2)$  are schema triples in  $G^\infty$ , then these two triples are also in  $(G_{/\equiv})^\infty$ , since  $S_G = S_{G_{/\equiv}}$  hence  $(S_G)^\infty = S_{G_{/\equiv}}^\infty$ .
- Else if both  $\phi(t_1)$  and  $\phi(t_2)$  are non-schema triples in  $G^\infty$ :
  - If both  $\phi(t_1)$  and  $\phi(t_2)$  are data triples in  $G^\infty$ , there exists a triple  $t'_1$  (resp  $t'_2$ ) in  $G$ , with same subject/object, from which  $\phi(t_1)$  (resp.  $\phi(t_2)$ ) is entailed using a sub-property constraint  $t_s^1$  (resp.  $t_s^2$ ) from  $S^\infty$ . Since  $\phi(t_1)$  and  $t'_1$  (resp.  $\phi(t_2)$  and  $t'_2$ ) have the same subject and object values, then  $t'_1$  and  $t'_2$  have same values on the places where  $t_1$  and  $t_2$  join. Therefore, if we assume that  $t'_1 = s_1 p_1 o_1$  and  $t'_2 = s_2 p_2 o_2$ , the  $G_{/\equiv}$  triples  $f(s_1) p_1 f(o_1)$  and  $f(s_2) p_2 f(o_2)$  necessarily have same values on the places where  $t_1$  and  $t_2$  join. Moreover, since  $S_G = S_{G_{/\equiv}}$  hence  $(S_G)^\infty = (S_{G_{/\equiv}})^\infty$ , these two  $G_{/\equiv}$  triples and the above-mentioned  $t_s^1$  and  $t_s^2$  schema triples, produce the counterpart triples of  $\phi(t_1)$  and  $\phi(t_2)$  in  $(G_{/\equiv})^\infty$ , which have same subject and object values. Thus, the  $q$  triples  $t_1$  and  $t_2$  embed in these two  $(G_{/\equiv})^\infty$  triples, if they embed in  $\phi(t_1)$  and  $\phi(t_2)$  in  $G^\infty$ .
  - Else if both  $\phi(t_1)$  and  $\phi(t_2)$  are type triples in  $G^\infty$ , say  $\phi(t_1) = u \tau c_1$  and  $\phi(t_2) = u \tau c_2$ , then  $t_1 = x \tau c_1$  and  $t_2 = x \tau c_2$  by definition of an RBGP query. As in the cases of single  $q$  triple embeddings,  $\phi(t_1)$  (resp.  $\phi(t_2)$ ) results either from a  $G$  triple  $u \tau c$  and a  $S_G^\infty$  triple  $c \prec_{sc} c_1$ , or a  $G$  triple  $u p u_1$  and a  $S_G^\infty$  triple  $p \leftrightarrow_d c_1$ , or a  $G$  triple  $u_1 p u$  and a  $(S_G)^\infty$  triple  $p \leftrightarrow_r c_1$ . Therefore, since  $S_G = S_{G_{/\equiv}}$  hence  $S_G^\infty = (S_{G_{/\equiv}})^\infty$ , for  $\phi(t_1)$  (resp.  $\phi(t_2)$ ), there are either a  $G_{/\equiv}$  triple  $f(u) \tau c$  and a  $(S_{G_{/\equiv}})^\infty$  triple  $c \prec_{sc} c_1$ , or a  $G_{/\equiv}$  triple  $f(u) p f(u_1)$  and a  $(S_{G_{/\equiv}})^\infty$  triple  $p \leftrightarrow_d c_1$ , or a  $G_{/\equiv}$  triple  $f(u_1) p f(u)$  and a  $(S_{G_{/\equiv}})^\infty$  triple  $p \leftrightarrow_r c_1$ , which entail  $f(u) \tau c_1$  and  $f(u) \tau c_2$  in  $G_{/\equiv}^\infty$ . Thus, the  $q$  triples  $t_1$  and  $t_2$  embed in these two  $G_{/\equiv}^\infty$  triples, if they embed in  $\phi(t_1)$  and  $\phi(t_2)$  in  $G^\infty$ .
  - Otherwise,  $\phi(t_1)$  is a data triple in  $G^\infty$  and  $\phi(t_2)$  is a type triple in  $G^\infty$ . This case is very similar to the two above case, hence we do not detail it.

- Otherwise,  $\phi(t_1)$  is a schema triple in  $G^\infty$  and  $\phi(t_2)$  is not a schema triples in  $G^\infty$ . In this case, since  $q$  is an RBGP query,  $q$  triples must be such that  $t_1$  is  $s_1 p_1 o_1$  with  $p_1 \in \{\prec_{sc}, \prec_{sp}, \leftrightarrow_d, \leftrightarrow_r\}$  and  $t_2$  is either  $s_2 p o_2$  or  $s_2 \tau c$ .

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

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

□

**Property 2.** *RBGP representativeness entails RBGP\* representativeness.*

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

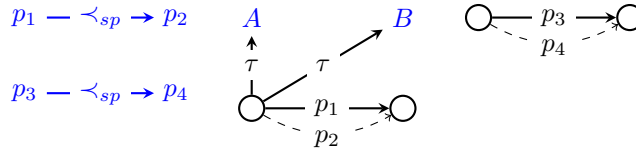
One can prove in a very similar fashion that summaries are also representative with respect to *regular path expression queries* such as SPARQL 1.1 property paths [33]: any regular path expression query whose results are non empty on an RDF graph  $G$  also has non-empty results on a summary of  $G$ . As a consequence, our summaries also preserve reachability in  $G$ : if a path leads from  $n_1$  to  $n_2$  in  $G$ , the same path also leads from the node representing  $n_1$  in the summary, to the node representing  $n_2$ .

We identify the following strong connection between two RDF graphs:

**Definition 8.** (STRONG ISOMORPHISM  $\simeq$ ) *A strong isomorphism between two RDF graphs  $G_1, G_2$  is an isomorphism which is the identity for the class and property nodes. We write this:  $G_1 \simeq G_2$ .*

From the definition of an RDF summary, it follows that any summary is (also) a summary of itself, more precisely  $(G_{/\equiv})_{/\equiv} \simeq G_{/\equiv}$ . We call this the **fixpoint property**; it captures the intuition that a summary cannot be summarized any further. The reason why  $(G_{/\equiv})_{/\equiv}$  and  $G_{/\equiv}$  are  $\simeq$  and not *identical* is that summarization creates “new” summary nodes (with fresh URIs). The summary preserves (reflects) the *structure* of the original graph, but it does not preserve the URIs of data nodes.

**Accuracy.** From the fixpoint property, it follows easily that any RDF summary is accurate (Definition 2). Indeed, for any graph  $G$  and RDF summary  $G_{/\equiv}$  of  $G$ , since  $(G_{/\equiv})_{/\equiv}$  and  $G_{/\equiv}$  are isomorphic, any query non-empty on  $(G_{/\equiv})_{/\equiv}$  is also non-empty on  $G_{/\equiv}$ .

Figure 5: RDF summary through  $\equiv_{fb}$  of the graph in Figure 4.

### 4.3 Bisimulation-based RDF summarization

We now instantiate our summarization framework. We start by defining *RDF FW (forward)*, *BW (backward)* and *FB bisimilarity*, based on Definition 6:

**Definition 9.** (FW RDF BISIMILARITY) *Forward (FW) RDF bisimilarity is an RDF node equivalence relation such that two data nodes  $n_1, n_2 \in \mathcal{G}$  are FW-bisimilar, denoted  $n_1 \equiv_{fw} n_2$ , iff either both  $n_1$  and  $n_2$  are not subjects of  $\mathcal{G}$  data triples or (i) for every  $\mathcal{G}$  data triple of the form  $n_1 \text{ p } m_1$  there exists a  $\mathcal{G}$  data triple of the form  $n_2 \text{ p } m_2$  such that  $m_1 \equiv_{fw} m_2$ , and conversely (ii) for every  $\mathcal{G}$  data triple of the form  $n_2 \text{ p } m_2$  there exists a  $\mathcal{G}$  data triple of the form  $n_1 \text{ p } m_1$  such that  $m_1 \equiv_{fw} m_2$ .*

Similarly, a *backward (BW) RDF bisimilarity*, noted  $\equiv_{bw}$ , is defined by inverting the triples' subject and object positions in the definition above. Further, a *forward and backward (FB) RDF bisimilarity* between data nodes, noted  $\equiv_{fb}$ , is defined by nodes that are both FW and BW similar. In Figure 4, each of the nodes  $A, B, p_1, p_2, p_3$  and  $p_4$  is  $\equiv_{fw}$ ,  $\equiv_{bw}$ , and  $\equiv_{fb}$  to itself (these are all RDF equivalence relations, thus they are reflexive), and none of them is bisimilar to any other node, since they are class and property nodes, and Definition 6 rules this out. Further,  $u_1$  is  $\equiv_{fw}$ ,  $\equiv_{bw}$ , and  $\equiv_{fb}$  to  $u_3$ , as in  $\mathcal{G}$ , they both have the properties  $p_1$  and  $\tau$  and no incoming properties; similarly  $u_2$  is equivalent to  $u_4$ . From the  $\equiv_{fw}$  perspective only,  $u_2, u_4$  and  $u_6$  are all equivalent as they have no outgoing edges. Bounded versions of these RDF bisimilarities could also be easily defined as discussed in Section 4.1 for their simple (non-RDF) counterparts.

So far, our equivalence notions do not exploit the possible types which may be attached to RDF nodes in a graph. To do so, we define:

**Definition 10.** (TYPED EQUIVALENCE) *Typed equivalence, noted  $\equiv_{\tau}$ , is an RDF node equivalence relation defined as follows: two data nodes  $n_1$  and  $n_2$  are type-equivalent, noted  $n_1 \equiv_{\tau} n_2$ , iff they have exactly the same set of types in  $\mathcal{G}$ , which is non-empty.*

In Figure 4,  $u_1 \not\equiv_{\tau} u_3$  as they have different types. A node in an RDF graph may have no type; such a node is equivalent to itself (again, due to the reflexivity of all RDF equivalence relations, including  $\equiv_{\tau}$ ), but not equivalent to any other node.

Clearly, each RDF node equivalence relation leads to a partition over the data nodes in  $\mathcal{G}$ . Thus, any RDF node equivalence (in particular,  $\equiv_{fw}$ ,  $\equiv_{bw}$ ,  $\equiv_{fb}$  or  $\equiv_{\tau}$  above) leads to a (possibly distinct) quotient summary, denoted  $\mathcal{G}/\equiv_{fw}$ ,  $\mathcal{G}/\equiv_{bw}$ ,  $\mathcal{G}/\equiv_{fb}$  and  $\mathcal{G}/\equiv_{\tau}$ . For illustration, Figure 5 shows the summary  $\mathcal{G}/\equiv_{fb}$  of the RDF graph  $\mathcal{G}$  in Figure 4. Since all nodes in that  $\mathcal{G}$  are untyped, its summary  $\mathcal{G}/\equiv_{\tau}$  has the same number of nodes as  $\mathcal{G}$ . This shows that summarization through  $\equiv_{\tau}$  alone may fail to reduce the graph size.

To exploit *type information when available, and structure similarity when types are lacking*, we define:

**Definition 11.** (TYPED BISIMILARITY SUMMARIES) *Let  $\equiv_{u, fw}$  (untyped FW equivalence) be an RDF node equivalence relation that holds between two nodes  $n_1, n_2$  iff (i)  $n_1, n_2$  have no types in  $\mathcal{G}$  and (ii)  $n_1 \equiv_{fw} n_2$ . The typed FW summary of  $\mathcal{G}$ , denoted  $\mathcal{G}/\equiv_{\tau, fw}$  is defined as:  $(\mathcal{G}/\equiv_{\tau})/\equiv_{u, fw}$ . The typed BW summary, noted  $\mathcal{G}/\equiv_{\tau, bw}$ , is similarly defined by replacing *fw* with *bw* in the above, while for the typed FB summary, noted  $\mathcal{G}/\equiv_{\tau, fb}$ , *fw* is replaced with *fb*.*



Let us follow for instance the construction of the typed FW summary. First, during summarization through  $\equiv_{\tau}$ , typed nodes from  $G$  are grouped according to their types, while untyped  $G$  nodes and edges between them are simply copied into an isomorphic summary subgraph. Then, we summarize again, applying  $\equiv_{fw}$  only to the untyped nodes of  $G_{/\equiv_{\tau}}$ . Overall, the process gives *priority to types*, i.e.,  $G$  nodes are summarized into a same node as soon as they have the same types in  $G$ , regardless of the surrounding graph structure; and second, for  $G$  nodes lacking types, we use structural bisimilarity to decide equivalence.

Figure 9 illustrates  $\equiv_{\tau,bw}$  summarization. As all data nodes are untyped in  $G$ ,  $u_1, u_4$  are summarized together in  $G_{/\equiv_{\tau,bw}}$ . In contrast, in  $G^{\infty}$ ,  $u_1$  is typed and  $u_4$  is not, therefore in  $(G^{\infty})_{/\equiv_{\tau,bw}}$  they are represented by different nodes.

## 5 Efficient summarization shortcuts

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

To avoid going through the saturation step, we identify a method (below) which guarantees that we can build RBGP\* representative summaries efficiently:

**Definition 12.** (SHORTCUT) *Summarization through the RDF node equivalence relation  $\equiv$  admits a shortcut iff for any RDF graph  $G$ ,  $(G^{\infty})_{/\equiv} \simeq ((G_{/\equiv})^{\infty})_{/\equiv}$  holds, where  $\simeq$  denotes a strong isomorphism (Definition 8).*

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

Figure 6 illustrates a shortcut for an  $\equiv_{fb}$ -based summary  $G_{/\equiv_{fb}}$  of an RDF graph  $G$ :  $(G^{\infty})_{/\equiv_{fb}} \simeq ((G_{/\equiv_{fb}})^{\infty})_{/\equiv_{fb}}$  holds. Note that the last summarization step is necessary here, because  $(G^{\infty})_{/\equiv_{fb}} \simeq (G_{/\equiv_{fb}})^{\infty}$  does not hold.

The next theorem is the major result of this work. It establishes a sufficient condition on an RDF node equivalence relation for the existence of a shortcut. To be able to state the theorem, we introduce the following:

**Representation function  $f$**  By the summary definition, to every node in  $G$  corresponds exactly one node in the summary  $G_{/\equiv}$ . We call *representation function* and denote  $f_{/\equiv}$  (or simply  $f$ , when this does not cause confusion) the function associating a summary node to each  $G$  node; we say  $f(n)$  *represents*  $n$  in the summary. For instance, for the graph  $G$  in Figure 4 and its summary in Figure 5,  $f(u_1) = f(u_3)$  is the node having the types  $A$  and  $B$  in Figure 5.

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

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

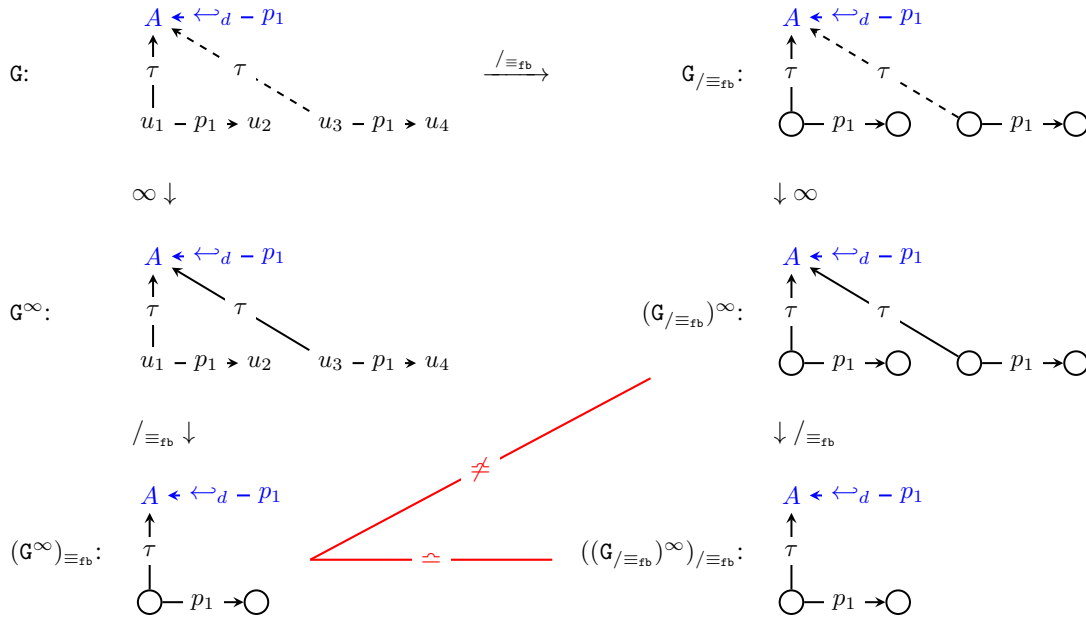


Figure 6: Shortcut example for RDF summarization through  $\equiv_{fb}$ .

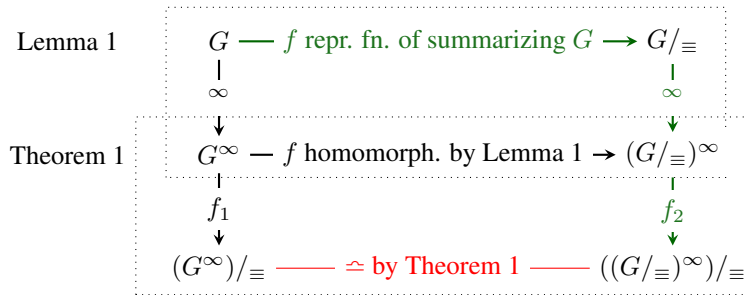


Figure 7: Illustration for Lemma 1 and Theorem 1.

The proof relies on the following lemma, which states a crucial structural property relating  $G$ ,  $G^\infty$  and the representation function  $f$  (see the upper part of Figure 7):

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

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

Recall from Section 2 that RDF saturation with RDFS constraints only adds edges between graph nodes, but does not add nodes. Thus, a node  $n$  is in  $G^\infty$  iff  $n$  is in  $G$ . Further, by the definition of our quotient-based summaries (Definition 7),  $n$  is in  $G$  iff  $f(n)$  is in  $G_{/\equiv}$ . Finally, again by the definition of saturation,  $f(n)$  is in  $G_{/\equiv}$  iff  $f(n)$  is in  $(G_{/\equiv})^\infty$ .

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

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

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

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

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

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

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

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

Lemma 1 is important as it connects RDF saturation (vertical  $\infty$ -labeled edges in Figure 7) to the representation of RDF nodes by corresponding summary nodes (black horizontal edges). The lemma states that the representation function associated to the summarization of  $G$  defines an homomorphism from the saturation of  $G$  to that of its summary. Note how this homomorphism (central horizontal edge in Figure 7) closes a rectangle materializing two paths from  $G$  to  $(G^\infty)_{/\equiv}$ : one consists of the first two steps of the shortcut (Definition 12), while the other requires saturating  $G$ .

*Proof.* (Theorem 1) We start by introducing some **notations** (see Figure 7). Let  $f_1$  be the representation function from  $G^\infty$  into  $(G^\infty)_{/\equiv}$ , and  $f_2$  be the representation function from  $(G_{/\equiv})^\infty$  into  $((G_{/\equiv})^\infty)_{/\equiv}$ .

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

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

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

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

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

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

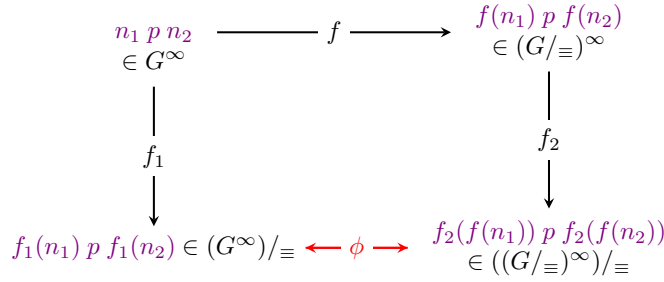


Figure 8: Diagram for the proof of Theorem 1.

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

For every edge  $n'_1 p n'_2$  in  $(G^\infty)_{/\equiv}$ , by definition of an RDF summary, there exists an edge  $n_1 p n_2$  in  $G^\infty$  such that  $n'_1 p n'_2 = f_1(n_1) p f_1(n_2)$ . Figure 8 illustrates the discussion. Further, if  $n_1 p n_2$  is in  $G^\infty$ , then  $f(n_1) p f(n_2)$  is in  $(G/\equiv)^\infty$  (Proposition 1), hence  $f_2(f(n_1)) p f_2(f(n_2))$  is in  $((G/\equiv)^\infty)_{/\equiv}$ . Therefore,

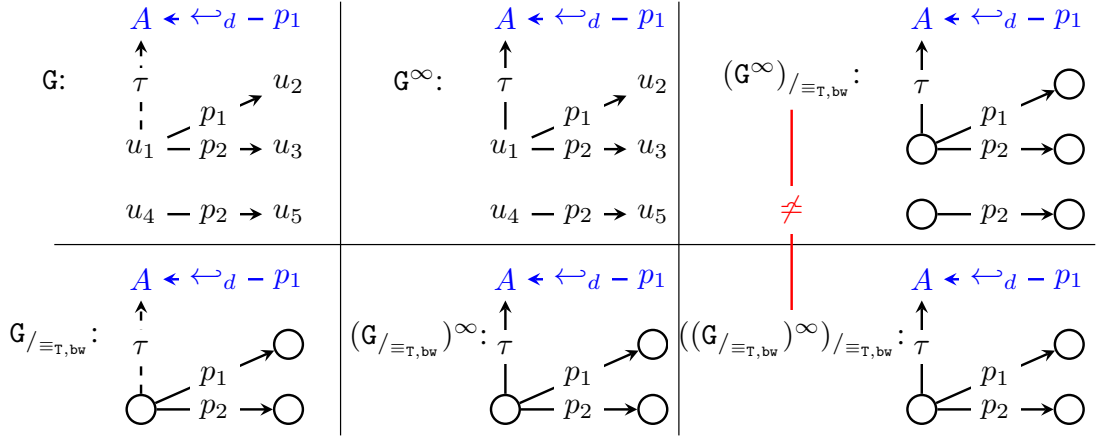
- since for every  $f_1(n_1) p f_1(n_2)$  edge in  $(G^\infty)_{/\equiv}$ , there is an edge  $f_2(f(n_1)) p f_2(f(n_2))$  in  $((G/\equiv)^\infty)_{/\equiv}$ , and
- since  $\varphi(f_1(n)) = f_2(f(n))$ , for  $n$  any  $G^\infty$  node, is a bijective function from all  $(G^\infty)_{/\equiv}$  nodes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

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

Similarly, for every edge  $n'_1 \tau c$  in  $(G^\infty)_{/\equiv}$ , by definition of an RDF summary, there exists an edge  $n_1 \tau c$  in  $G^\infty$  such that  $n'_1 \tau c = f_1(n_1) \tau c$ . Further, if  $n_1 \tau c$  is in  $G^\infty$ , then  $f(n_1) \tau c$  is in  $(G/\equiv)^\infty$  (Proposition 1), hence  $f_2(f(n_1)) \tau c$  is in  $((G/\equiv)^\infty)_{/\equiv}$ . Therefore,

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

Now, for every edge  $n''_1 \tau c$  in  $((G/\equiv)^\infty)_{/\equiv}$ , by definition of an RDF summary, there exists an edge  $n'_1 \tau c$  in  $(G/\equiv)^\infty$  such that  $n''_1 \tau c = f_2(n'_1) \tau c$ . Hence, by Proposition 1, there exists an edge  $n_1 \tau c$  in  $G^\infty$  such that  $n'_1 \tau c = f(n_1) \tau c$ . Moreover, since  $n_1 \tau c$  is in  $G^\infty$ ,  $f_1(n_1) \tau c$  is in  $(G^\infty)_{/\equiv}$ . Therefore, since for every  $f_2(f(n_1)) \tau c$  edge in  $((G/\equiv)^\infty)_{/\equiv}$ , there is an edge  $f_1(n_1) \tau c$  in  $(G^\infty)_{/\equiv}$ ,

Figure 9: Sample typed bisimilarity summarization through  $\equiv_{T,bw}$ .

and since  $\varphi(f_1(n)) = f_2(f(n))$ , for  $n$  any  $G^\infty$  node, is a bijective function from all  $(G^\infty)_{/≡}$  nodes to all  $((G_{/≡})^\infty)_{/≡}$  nodes,  $(G^\infty)_{/≡}$  contains the image of all  $((G_{/≡})^\infty)_{/≡} n''_1 \tau c$  triples through  $\varphi^{-1}$  (\*\*').

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

Importantly, *not all RDF equivalence relations admit a shortcut*, as Figure 9 shows. From  $G$ , the top row traces the two-step computation of  $(G^\infty)_{/≡_{T,bw}}$ , while moving down to the bottom row we see the three-step shortcut to  $((G_{/≡_{T,bw}})^\infty)_{/≡_{T,bw}}$ . The shortcut does not lead to (a graph  $\simeq$  to) the  $\equiv_{T,bw}$  summary of  $G^\infty$ , because  $u_1$  and  $u_4$  are kept separate in the top row, while they are fused on the bottom one.

However, for RDF summaries based (only) on RDF bisimilarity relations, we readily establish based on Theorem 1:

**Theorem 2.** *Summarization through the RDF equivalence relations  $\equiv_{fw}$ ,  $\equiv_{bw}$  and  $\equiv_{fb}$  admits a shortcut.*

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

## 6 Data property cliques and clique-based equivalence

This section sets the first step for the novel contributions of this paper: the novel notion of property cliques (Section 6.1), followed by two new RDF node equivalence relations (Section 6.2) on which we build our summaries.

### 6.1 Data property cliques

The simplest relation between two data properties in a graph is *co-occurrence*, when a resource is the source and/or target of two data properties. However, in heterogeneous RDF graphs, two properties  $p_1, p_2$  may co-occur on some nodes, while other nodes may have  $p_1$  or only  $p_2$ , while such nodes have yet other properties  $p_3, p_4$  etc. Thus, we introduce a more general and flexible notion crucial to this work, as follows:

**Definition 13.** (PROPERTY RELATIONS AND CLIQUES) *Let  $p_1, p_2$  be two data properties in  $G$ :*

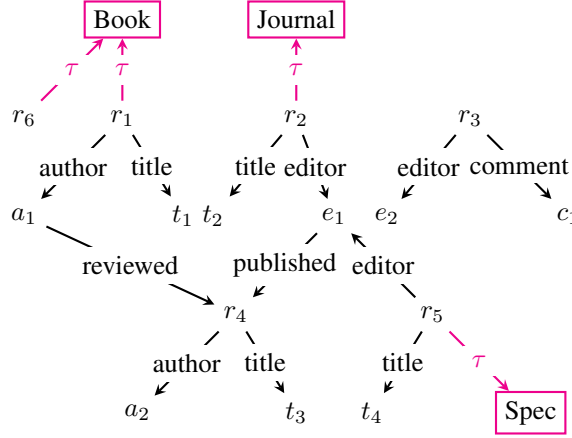


Figure 10: RDF graph used for clique-based summarization.

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

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

For illustration, consider the sample graph in Figure 10. Here, properties  $a$  and  $t$  are source-related due to  $r_1$  (condition (i) in the definition). Similarly,  $t$  and  $e$  are source-related due to  $r_2$ ; consequently,  $a$  and  $e$  are source-related (condition (ii)). Properties  $r$  and  $p$  are target-related due to  $r_4$ . The non-empty source cliques are  $SC_1 = \{a, t, e, c\}$ ,  $SC_2 = \{r\}$  and  $SC_3 = \{p\}$ , whereas the non-empty target cliques are  $TC_1 = \{a\}$ ,  $TC_2 = \{t\}$ ,  $TC_3 = \{e\}$ ,  $TC_4 = \{c\}$  and  $TC_5 = \{r, p\}$ .

A source clique can be seen as “all the data properties of same-kind resources”; above, it makes sense to group together properties corresponding to various kinds of publications, such as author, title, editor, and comment, *despite the heterogeneity exhibited by the presence/absence of these properties* on the RDF graph nodes. Similarly, a target clique comprises “the data properties of which same-kind resources are values”.

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

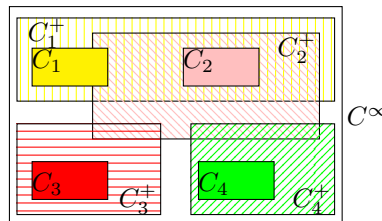
The target and source cliques of the resources in the graph shown in Figure 10 are shown in Table 1.

**Definition 14.** (PROPERTY DISTANCE IN A CLIQUE) *The distance between two data properties  $p, p'$  in a source (resp. target) clique is:*

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

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

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

Figure 11: Sample cliques of  $G$ , their saturations, and their enclosing clique in  $G^\infty$ .

In Figure 10, the distance between  $a$  and  $t$  is 0 since  $r_1$  has both. The distance between  $a$  and  $e$  is 1, while the distance between  $a$  and  $c$  is 2. Clearly,  $p$  and  $p'$  are at distance  $n$  for  $n > 0$  iff a resource has both  $p$  and  $p''$ , and further  $p''$  is at distance  $n - 1$  from  $p'$ .

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

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

**Lemma 2** (Saturation vs. property cliques). *Let  $C, C_1, C_2$  be distinct non-empty source (or target) cliques of  $G$ .*

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

$$D_1^+ \cap D_2^+ \neq \emptyset, \dots, D_{n-1}^+ \cap D_n^+ \neq \emptyset$$
4. *Let  $p_1, p_2$  be two data properties in  $G$ , whose source (or target) cliques are  $C_1$  and  $C_2$ . Properties  $p_1, p_2$  are in the same source (resp. target) clique  $C^\infty$  corresponding to  $G^\infty$  if and only if there exist  $k$  non-empty source (resp. target) cliques of  $G$ ,  $k \geq 0$ , denoted  $D_1, \dots, D_k$  such that:*



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

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

1. Any resource  $r \in G$  having two data properties also has them in  $G^\infty$ ; thus, any data properties in the same source clique in  $G$  are also in the same source clique in  $G^\infty$ . The unicity of  $C^\infty$  is ensured by the fact that the source cliques of  $G^\infty$  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  $G^\infty$ ,  $r_1$  has  $p_1$  and  $p$ , thus these two properties are in the same  $G^\infty$  clique. Similarly,  $r_2$  has  $p_2$  and  $p$ , which ensures that  $p$  is also in the same  $G^\infty$  source clique.
3. Let  $\{p_1, \dots, p_k\}$  be the data properties that appear both in  $G$  and in  $C^\infty$ ; it follows from the saturation rules and the definition of cliques, that  $k > 0$ . For  $1 \leq i \leq k$ , let  $C_i$  be the  $G$  source clique comprising  $p_i$ . Applying lemma point 1.,  $C_i \subseteq C^\infty$  for each  $1 \leq i \leq k$ . Further, it is easy to see that  $C_i^+ \subseteq C^\infty$ , since any property that saturation adds to  $C_i^+$  is also added by saturation to  $C^\infty$ . Thus,  $\bigcup_{1 \leq i \leq k} C_i^+ \subseteq C^\infty$ .

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

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

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

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

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

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

The lemma is crucial as it states the relationship between the cliques (thus, the summaries) of  $G$  and  $G^\infty$ ; our goal is to build the summary of the latter *without* having access to it, but just access to  $G$ . Figure 11 illustrates the lemma. The clique  $C^\infty$  of  $G^\infty$  encloses the cliques  $C_1, C_2, C_3, C_4$  of  $G$  (Lemma item 1).  $C_2^+$  intersects all of  $C_1^+, C_3^+$  and  $C_4^+$ , thus they are all in the same clique  $C^\infty$  (item 2), which is the union of  $C_1^+, C_2^+, C_3^+$  and  $C_4^+$  (item 3). A chain of intersecting cliques connects e.g.,  $C_1^+$  through  $C_2^+$  to  $C_4^+$ . Item 4 is not illustrated, as individual properties do not show; it follows quite directly from item 3.

## 6.2 Clique-based RDF node equivalences

Building on property cliques, we define the following equivalences:

**Definition 15.** (STRONG AND WEAK EQUIVALENCE) Strong (denoted  $\equiv_s$ ) and weak (denoted  $\equiv_w$ ) equivalences are RDF node equivalence relations. Two data nodes of  $G$  are strongly equivalent, denoted  $n_1 \equiv_s n_2$ , iff they have the same source and target cliques.

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

Observe that strong equivalence implies weak equivalence.

In Figure 10, the resources  $r_1, r_2, r_3, r_5$  are *strongly* equivalent to each other, as well as  $t_1, t_2, t_3, t_4$ . Moreover,  $r_1, \dots, r_5$  are *weakly* equivalent to each other due to their common source clique  $SC_1$ , as well as  $t_1, t_2, t_3, t_4$  due to their common target clique; the same holds for  $a_1$  and  $a_2$ , and separately for  $e_1$  and  $e_2$ . In general, weakly equivalent resources can be connected as exemplified in Figure 12, by an *alternating sequence of source and target cliques*.

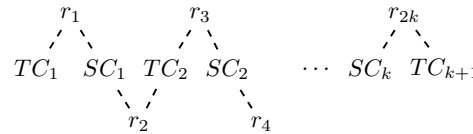


Figure 12: Weakly equivalent resources of  $G$  and their cliques.

## 7 Weak-equivalence summaries

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

### 7.1 Weak summary

Our first new summary is solely based on weak equivalence:

**Definition 16.** (WEAK SUMMARY) The weak summary of the graph  $G$ , denoted  $G/w$ , is its quotient graph w.r.t. the weak equivalence relation  $\equiv_w$ .

For each set of  $\equiv_w$ -equivalent  $G$  nodes, there is exactly one node in  $G/w$ . Note that the partition of  $G$  nodes into  $\equiv_w$  equivalence classes is also a *partition of  $G$  data properties at the source*, that is: the sources of  $G$  edges labeled with a given data property  $p$  are all weakly equivalent. For instance, in Figure 10, the sources of all “editor” edges are in the weak equivalence class  $\{r_1, r_2, r_3, r_4, r_5\}$ .

The  $\equiv_w$  partition over the set of nodes of  $G$  also induces a *partition of the source cliques of  $G$* . Indeed, if  $n_1 \equiv_w n_2$ , their source cliques are connected through a chain of alternating cliques as sketched in Figure 12; conversely, such a chain, by definition, only connects weakly related resources. By a symmetrical reasoning, the weak equivalence classes of  $G$  also lead to a *partition over the target cliques of  $G$* . For instance, in Figure 10, to the equivalent resource set  $\{r_1, \dots, r_5\}$  corresponds the set of source cliques  $\{SC_1\}$  and the set of target cliques  $\{TC_5\}$ . Thus, *to each set  $S$  of weakly equivalent  $G$  nodes one can associate through a bijection, a set of  $G$  source cliques, and a set of  $G$  target cliques*. This leads to:

**Definition 17.** (SET REPRESENTATIVE) Let  $N$  be any injective function taking as input two sets of URIs, and returning a new URI. Let  $S$  be an equivalence class of  $G$  data nodes w.r.t.  $\equiv_w$ . The weak summary node representing all  $S$  nodes, also called the set representative of  $S$  in the summary, is:  $N(\bigcup_{r \in S} TC(r), \bigcup_{r \in S} SC(r))$ .

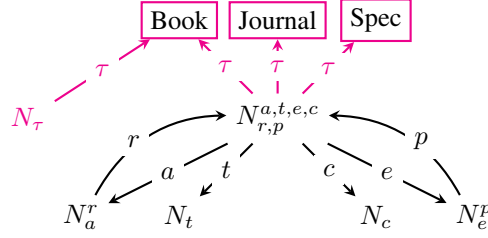


Figure 13: Weak summary of the RDF graph in Figure 10.

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

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

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

In the particular case where a data resource  $r \in G$  is neither the source nor the target of data properties, i.e.,  $TC(r) = SC(r) = \emptyset$  (thus  $r$  can only appear in  $\tau$  triples),  $r$  is represented by  $N(\emptyset, \emptyset)$  which we denote  $N_\tau$  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  $N_\tau$ .

The weak summary of the graph in Figure 10 is shown in Figure 13. Its nodes are:  $N_{r,p}^{a,t,e,c}$  for the relatedness partition set  $\{r_1, \dots, r_5\}$ . The target properties of this node are  $TC(r_4)$  since the other nodes have empty target clique; the source properties are those in  $SC(r_1)$  which is also the source clique of all the other resources in the set;  $N_a^r$  for the set  $\{a_1, a_2\}$ ;  $N_t$  for the relatedness partition set  $\{t_1, t_2, t_3, t_4\}$  etc. The edges from  $N_{r,p}^{a,t,e,c}$  to  $N_a^r$  and  $N_t$  copy the outgoing edges of  $r_1$ , represented by  $N_{r,p}^{a,t,e,c}$ ; the edge to  $N_e^p$  is due to  $r_2$  and  $r_3$ ; the edge to  $N_c$  is due to  $r_3$ . The edge from  $N_a^r$  to  $N_{r,p}^{a,t,e,c}$  is due to  $a_1$ , and the edge from  $N_e^p$  to  $N_{r,p}^{a,t,e,c}$  is due to  $e_1$ . The  $\tau$  edges outgoing  $N_{r,p}^{a,t,e,c}$  are due to the resources  $r_1, r_2$  and  $r_3$ ; the creation of  $N_\tau$  (shown in purple font) is due to the node  $r_6$  in the original graph.

The weak summary has the following important properties:

**Property 3.** (UNIQUE DATA PROPERTIES) *Each  $G$  data property appears exactly once in  $G_w$ .*

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

Importantly, the above Property 3 warrants that the number of data edges in  $G_w$  is exactly  $|D_G|_p^0$ , the number of distinct data properties in  $G$ . Thus, its number of data nodes is at most  $2|D_G|_p^0$ . The number of type triples in  $G_w$  is bound by  $\min(|T_G|_e, 2|D_G|_p^0 * |T_G|_o^0)$ : the latter corresponds to the case when every data node in  $G_w$  is of every type in  $T_G$ .

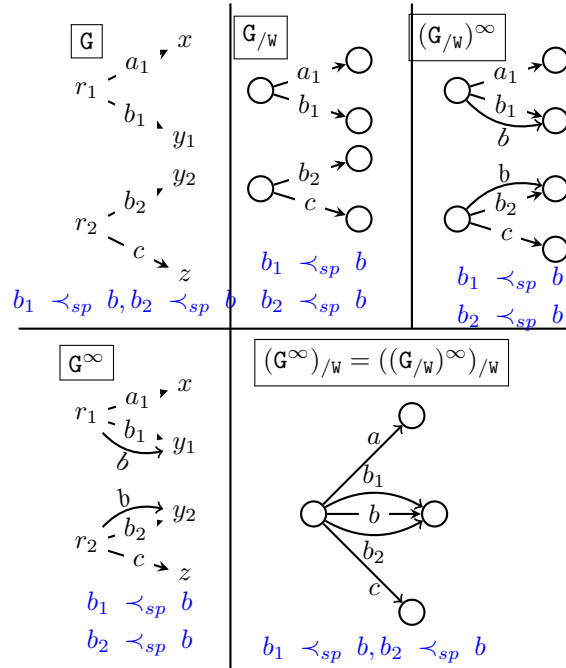


Figure 14: Weak summary shortcut illustration (summary nodes shown as unlabeled circles).

The next important property enjoyed by weak summaries is the existence of *shortcuts*, in the sense of Definition 12. Figure 14 exemplifies this, by tracing the transformation of a graph  $G$  on one hand, into  $G/w$ ,  $(G/w)^\infty$ , and then  $((G/w)^\infty)/w$ ; and on the other hand, from  $G$  to  $G^\infty$  then  $(G^\infty)/w$ . The graph at the bottom right in the figure is at the same time  $((G/w)^\infty)/w$  and  $(G^\infty)/w$ .

Showing that  $G/w$  admits a shortcut is rather involved. For this, we rely on the next Proposition:

**Property 4.** (SAME CLIQUES-W)  $G^\infty$  and  $(G/w)^\infty$  have identical source clique sets, and identical target clique sets. Further, any node  $n$  in  $G^\infty$  has the source clique  $SC$  and target clique  $TC$  iff  $f(n)$  in  $(G/w)^\infty$  has the source clique  $SC$  and target clique  $TC$ , with  $f$  the representation function (recall Proposition 1).

Proposition 4 is established based on Proposition 1, and on the following Lemma:

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

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

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

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

- For every  $1 \leq j \leq m$ :
  - $TC_W^j$  is the target clique of a  $G/w$  node  $n^j$ ;
  - $n^j$  represents a set of weakly-equivalent  $G$  resources, which are targets only of properties in  $TC_W^j$ . Thus, the properties in  $TC_W^j$  are target-related in  $G$ .

- Thus, in  $G^\infty$ , also, the properties in  $TC_W^j$  are target-related.
- From this and the definition of a saturated graph and of a saturated target clique, it follows that the properties from  $(TC_W^j)^+$  are target-related in  $G^\infty$ .
- Further, still by Lemma 2, point 4, each  $(TC_W^j)^+$  intersects at least another  $(TC_W^l)^+$  for  $1 \leq l \neq j < m$ , thus the target properties in all the  $(TC_W^j)^+$  for  $1 \leq j \leq m$ , and in particular  $p$ , are target-related to each other in  $G^\infty$ . Thus,  $p$  is target-related in  $G^\infty$  to all properties from  $TC_W^\infty$ .

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

The existence of a shortcut for  $W$ -summaries, stated by Theorem 3 below, is then shown based on Proposition 4 by proving that  $W$ -summaries enjoy the sufficient condition stated by Theorem 1.

**Theorem 3** ( $W$  summary shortcut). *Weak summaries admit a shortcut.*

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

Two summary nodes are weakly-equivalent if they have the same non-empty source or target clique. The claim to prove hence immediately follows from Property 4, which guarantees that  $n_1$  and  $f(n_1)$ , resp.  $n_2$  and  $f(n_2)$ , have the same source and target cliques.  $\square$

**Weak summary size.** The size of the weak summary is at most that of the original graph. This bound is reached if all nodes have distinct source and target properties. Moreover, Property 3 entails that the number of data edges in  $G/W$  is exactly the number of distinct data properties in  $G$ . Thus, its number of data nodes is at most twice this number. The number of type triples in  $G/W$  is bound by the smallest between the number of type edges in  $G$ .

## 7.2 Typed weak summary

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

**Definition 18.** (TYPED EQUIVALENCE) Typed equivalence, noted  $\equiv_T$ , is an RDF node equivalence relation defined as follows: two data nodes  $n_1$  and  $n_2$  are type-equivalent, noted  $n_1 \equiv_T n_2$ , iff they have exactly the same set of types in  $G$ , which is non-empty.

**Definition 19.** (TYPED WEAK SUMMARY) Let  $\equiv_{UW}$  (untyped weak equivalence) be an RDF node equivalence relation that holds between two nodes  $n_1, n_2$  iff (i)  $n_1, n_2$  have no types in  $G$  and (ii)  $n_1 \equiv_w n_2$ . The typed weak summary  $G_{TW}$  of an RDF graph  $G$  is the untyped-weak summary of the type-based summary of  $G$ , namely  $(G/T)_{UW}$ .

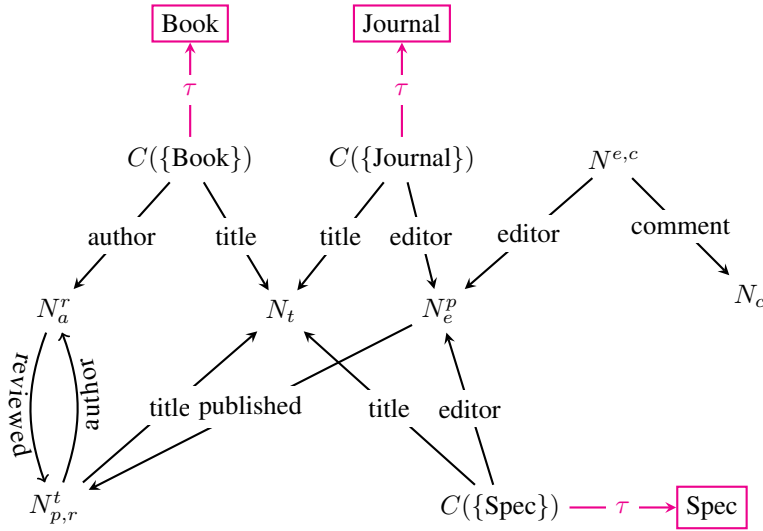


Figure 15: Typed weak summary  $G_{TW}$  of the graph  $G$  in Figure 10.

Untyped weak equivalence is the weak equivalence  $\equiv_w$  restricted to untyped resources only. The typed-weak summary  $G_{UW}$  treats the  $G$  triples whose subject and object lack types exactly like  $G_w$ , and groups typed resources by their types. For example, Figure 15 shows the typed weak summary of the RDF graph in Figure 10. It represents all editors with a single node, all titles with a single node etc., like  $G_w$  (Figure 13). However, in contrast with  $G_w$ ,  $G_{TW}$  represents publications by three distinct nodes, depending on their types.

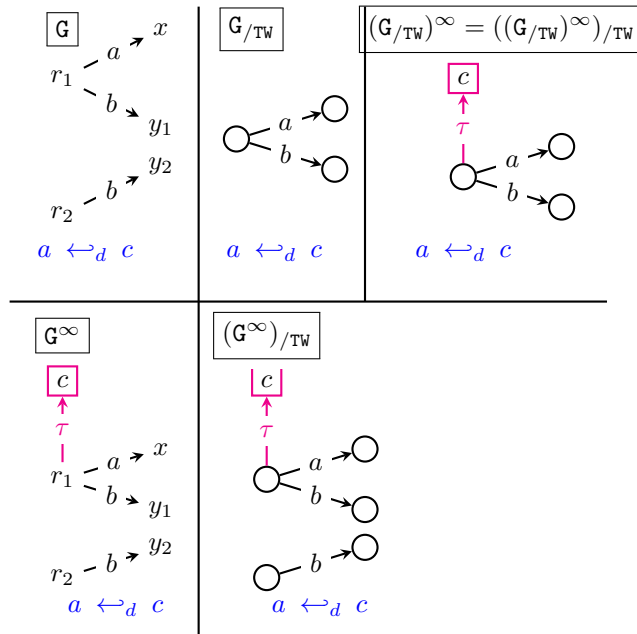


Figure 16: Typed weak summary shortcut counter-example.

**Property 5** (No shortcut for  $G_{TW}$ ). *Typed weak summaries do not admit a shortcut.*

Figure 16 shows a counter-example. In  $G$  and  $G_{TW}$ , resources are untyped; a typed resource only appears due to saturation and the constraint  $a \leftrightarrow_d c$ . Thus, in  $G_{TW}$ , one (untyped) node represents all data property subjects; in  $(G_{TW})^\infty$ , this (single) node gains a type, and in  $((G_{TW})^\infty)_{TW}$ , it is represented by a single node. In contrast, in  $G^\infty$ , one resource is typed and the other one isn't, leading to distinct nodes in  $(G^\infty)_{TW}$ . This cannot be isomorphic with the result obtained from  $G$  if one starts by TW summarization, thus by Definition 12,  $G_{TW}$  does not admit a shortcut.

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

## 8 Strong equivalence summaries

We now present summaries based on the strong RDF node equivalence. Strong summaries (Section 8.1) rely on it alone, while typed strong summaries (Section 8.2) give preeminence to types.

### 8.1 Strong summary

**Definition 20.** (STRONG SUMMARY) *The strong summary of the graph  $G$ , denoted  $G_S$ , is its quotient graph w.r.t. the strong equivalence relation  $\equiv_S$ .*

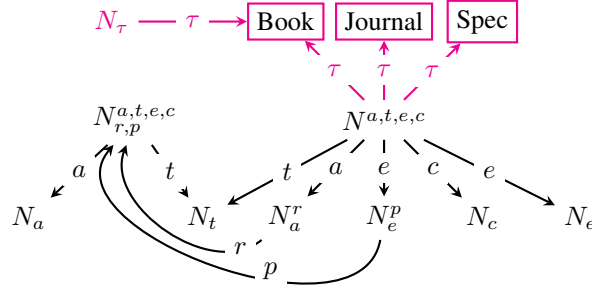


Figure 17: Strong summary of the RDF graph in Figure 10.

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 10 is shown in Figure 17. The strong summary comprises the nodes:

- $N(\emptyset, SC_1) = N^{a,t,e,c}$ , for  $r_1, r_2, r_3$  and  $r_5$ ;
- $N(TC_5, SC_1) = N^{a,t,e,c}_{r,p}$ , for  $r_4$ ;
- $N(TC_1, SC_2) = N_r^a$ , for  $a_1$ ;
- $N(TC_2, \emptyset) = N_t$ , for  $t_1, t_2, t_3$  and  $t_4$ ;

- $N(TC_3, SC_3) = N_e^p$ , for  $e_1$ ;
- $N(TC_3, \emptyset) = N_e$ , for  $e_2$ ;
- $N(TC_4, \emptyset) = N_c$ , for  $c_1$ ;
- $N(TC_1, \emptyset) = N_a$ , for  $a_2$ ;
- $N(\emptyset, \emptyset) = N_\tau$  for  $r_6$ ;
- Book, Journal, and Spec are copied from the schema of  $G$ .

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

Similarly to weak summaries, strong summaries admit a shortcut (Definition 12); this is our last major formal result. We prove this based on counterparts of statements established for  $G_W$ :

**Property 6.** (SAME CLIQUES-S)  $G^\infty$  and  $(G/S)^\infty$  have identical source clique sets, and identical target clique sets. Further, any node  $n$  in  $G^\infty$  has the source clique  $SC$  and target clique  $TC$  iff  $f(n)$  in  $(G/S)^\infty$  has the source clique  $SC$  and target clique  $TC$ , with  $f$  the representation function (recall Proposition 1).

The above Proposition follows directly from Proposition 1 and the following Lemma:

**Lemma 4.** Data properties are target-related (resp. source-related) in  $(G/S)^\infty$  iff they are target-related (resp. source-related) in  $G^\infty$ .

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

“Only If”: if two data properties  $p_1$  and  $p_2$  are target-related in  $(G/S)^\infty$ , then they belong to a same target clique  $TC^{S,\infty}$ , in which they are at distance  $n \geq 0$ , i.e., they are target-related because of a set  $\bigcup_{i=0}^n \{r_{i+1}\}$  of nodes (recall Definition 14) which all have the target clique  $TC^{S,\infty}$ . In  $G/S$ , each such  $r_{i+1}$  has a target clique  $TC_i^S \subseteq TC^{S,\infty}$ , moreover each  $r_{i+1}$  results from a set of  $G$  nodes  $n_{i+1}^j$ ,  $j \geq 1$ , which by definition of a strong RDF summary, have all the source clique  $TC_i^S$ . Hence, every such  $n_{i+1}^j$  node has target clique  $TC^{S,\infty}$  in  $G^\infty$  (since  $G$  and  $G/S$  have the same schema), in which  $p_1$  and  $p_2$  are target related.  $\square$

Based on Proposition 6, we prove that  $G_S$  enjoys the sufficient condition stated by Theorem 1:

**Theorem 4** (Strong summary shortcut). *Strong summaries admit a shortcut.*

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

Two summary nodes are strongly-equivalent if they have the same non-empty source or target clique. The claim to prove hence immediately follows from Property 4, which guarantees that  $n_1$  and  $f(n_1)$ , resp.  $n_2$  and  $f(n_2)$ , have the same source and target cliques.  $\square$



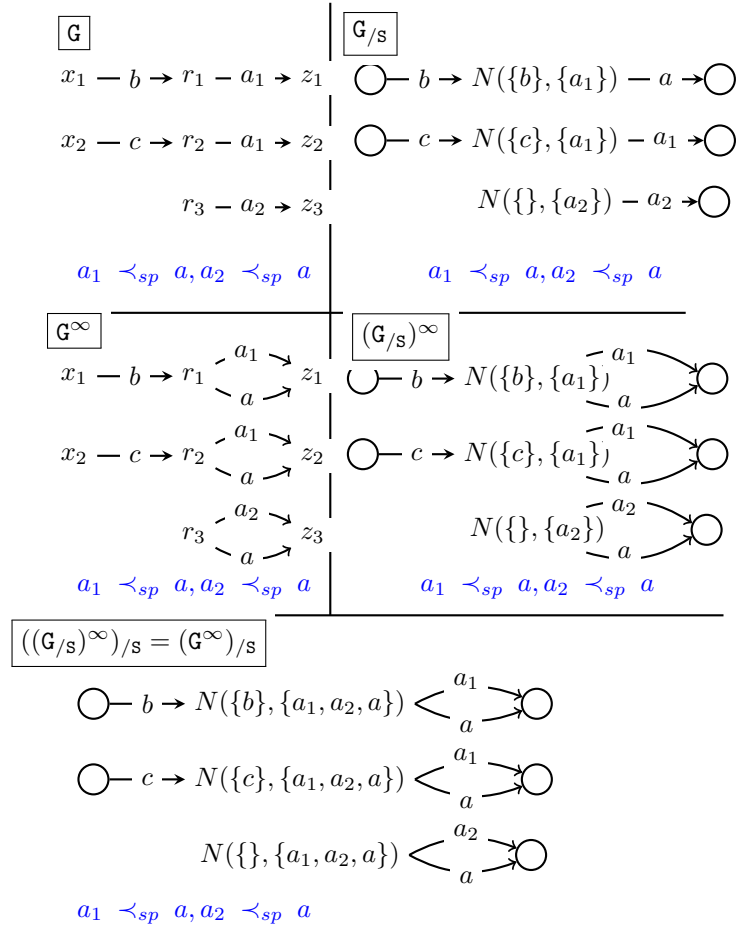


Figure 18: Illustration of the shortcut on a strong summary (some summary nodes are shown as unlabeled circles).

Figure 18 illustrates this on an example.

**Strong summary size.** 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). 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  $S_G$ .

## 8.2 Typed strong summary

The summary presented here is the counterpart of the typed weak one from Section 7.2, but based on strong equivalence.

**Definition 21.** (TYPED STRONG SUMMARY) *Let  $\equiv_{US}$  (untyped strong equivalence) be an RDF node equivalence relation that holds between two nodes  $n_1, n_2$  iff (i)  $n_1, n_2$  have no types in  $G$  and (ii)  $n_1 \equiv_S n_2$ . The typed strong summary  $G_{/TS}$  of an RDF graph  $G$  is defined as:  $(G_{/T})_{/US}$ .*

Untyped strong equivalence restricts strong equivalence to untyped resources only. The summary  $G_{/TS}$  summarizes untyped data resources in strong fashion, and groups all typed resources by their types. In our example, it turns out that  $G_{/TS}$  for the RDF graph in Figure 10 coincides with the  $G_{/TW}$  one shown in Figure 15. As can be easily seen from their definitions, the type-weak and type-strong summaries behave identically on the triples involving typed resources; on the untyped ones, the difference is of the same nature as the difference between the strong and weak summaries.

**Property 7** (No shortcut for  $G_{/TS}$ ). *Typed strong summaries do not admit a shortcut.*

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

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

## 9 Summarization algorithms

We have implemented a summarization tool in Java 1.8 (approx. 34.000 lines) which, given a graph  $G$ , builds  $G_{/W}$ ,  $G_{/TW}$ ,  $G_{/S}$  and  $G_{/TS}$  as well as the bisimulation-based RDF summaries  $G_{/fw}$ ,  $G_{/bw}$  and  $G_{/fb}$ .

### 9.1 Bisimulation summaries

We follow the algorithm from [19] for summarizing the data triples, and add our specific treatment of schema and type triples to obtain representative summaries.  $G$  is initially stored in PostresSQL, in a dictionary-encoded triple table. The bisimulation summary itself is stored in MapDB, allowing to build the summary in memory, backed by disk storage. Inspired by [23], we implement the signature table as a Trie, which we choose to keep in memory.

Algorithm 1 shows how to generate the summary  $G_{\equiv_{fb}}$  of an RDF graph  $G$ .  $N_G^D$  and  $E_G$  denote the data nodes of  $G$  and the edges of  $G$ , while the data nodes, respectively edges of the summary are  $N_{fb}^D$  and  $E_{fb}$ . The algorithm is iterative, first representing all  $G$  data nodes by a single summary node, and gradually replacing this initial representation function  $f_0$  by  $f_1, f_2$  etc. until stationarity is reached;  $f_k(n)$  is the representation function of a node  $n$  in round  $k$ . In each round  $k$ , we compute the  $k$ -th signature of

**Input** : RDF graph  $G = \langle D_G, S_G, T_G \rangle$   
**Output**: Forward and backward bisimulation summary  $G_{\equiv_{fb}}$   
 $f_0(n) \leftarrow 0$  for  $\forall n \in N_G^D$ ;  $seed \leftarrow 1$ ;  $sigTable \leftarrow \{\}$ ;  $k \leftarrow 0$ ;  
**repeat**  
   $k \leftarrow k + 1$   
  **for**  $n \in N_G^D$  **do**  
     $f_k(n) \leftarrow \text{getRepresentative}(n, D_G, k)$   
  **end**  
**until**  $f_k(n) = f_{k-1}(n)$  for all  $n \in N_G^D$ ;  
 $N_{fb}^D \leftarrow \{f_k(n) \mid n \in N_G^D\}$ ;  $E_{fb} \leftarrow \{(f_k(n), f_k(n')) \mid (n, n') \in E_G\}$ ;  
**Function**  $\text{getRepresentative}(n, D_G, k)$   
   $FW_k(n) = \{(p, f_{k-1}(o)) \mid n p o \in D_G\}$ ,  $BW_k(n) = \{(p, f_{k-1}(s)) \mid s p n \in D_G\}$   
  Sort  $FW_n$  and  $BW_n$  by their first, then by their second component  
   $sig_k(n) = f_{k-1}(n) \| FW_k(n) \| BW_k(n)$  //  $\|$  denotes concatenation  
   $repr \leftarrow sigTable.get(sig_k(n))$  // look up representative of  $n$  in round  $k$   
  **if**  $repr = \perp$  **then**  
     $repr \leftarrow seed$ ;  $sigTable.put(sig_k(n), repr)$ ;  $seed \leftarrow seed + 1$   
  **end**  
  return  $repr$ ;  
**endFunction**

**Algorithm 1:** Bisimulation summarization algorithm for  $G_{\equiv_{fb}}$  based on [19].

a data node  $n$  from based on data triples of which  $n$  is a subject and/or object, and represent  $n$  in round  $k$  by a summary node corresponding to this signature. Thus, there is a forward component ( $FW_k(n)$ ) and a backward component ( $BW_k(n)$ ) of the signature, each of which also reflects the current representation of the node at the other end of the triple. The  $k$ -th signature of  $n$  concatenates  $f_{k-1}(n)$ ,  $FW_k(n)$  and  $BW_k(n)$ . To build  $G_{\equiv_{fw}}$  or  $G_{\equiv_{bw}}$ , the signature should use only  $FW_n$ , or only  $BW_n$ , respectively.

## 9.2 Clique summaries

The clique summaries are built mostly in memory and also using a few SQL queries. The algorithms for building  $G_w$  and  $G_{TW}$  are significantly different, and more efficient, than what could be designed by simply following the summary definitions; we outline them below.  $G_{/S}$  and  $G_{/TS}$  are built in a fashion closer to the definitions. Below, we first introduce the data structures used by the algorithms.

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

Figure 19: Overview of the data structures for different summaries.

```

Input : Data triples table  $D_G$ , summary  $G_{/W}$ 
Output: Data triples represented in  $G_{/W}$ 
foreach  $s p o$  in  $D_G$  do
   $src \leftarrow \text{getSource}(s, p, G_{/W}); \text{targ} \leftarrow \text{getTarget}(o, p, G_{/W})$ 
  //  $\text{getTarget}$  may have modified  $src$  and vice-versa
   $src \leftarrow \text{getSource}(s, p, G_{/W}); \text{targ} \leftarrow \text{getTarget}(o, p, G_{/W})$ 
  if  $\text{!existsDataTriple}(G_{/W}, src, p, targ)$  then
    |  $\text{createDataTriple}(G_{/W}, src, p, targ)$ 
  end
end
Procedure  $\text{createDataTriple}(G_{/W}, src, p, targ)$ 
  |  $\text{dtp.put}(p, src \ p \ targ)$ 
  |  $\text{dpSrc.put}(p, src); \text{srcDps.put}(src, p)$ 
  |  $\text{dpTarg.put}(p, targ); \text{targDps.put}(targ, p)$ 
endProcedure

```

**Algorithm 2:** Summarizing data triples in  $G_{/W}$ 

### 9.3 Data structures

Currently, our summary graphs are built in memory, based on the Trove library<sup>4</sup> providing efficient collection data structures. We represent  $G$  URIs by integers, and store them in a set of map/ multi-map data structures, listed in Figure 19, together with the summaries which use them. All structures are not needed by all summaries, but many of them are used repeatedly. We describe them below, underlined for visibility.

rd maps  $G$  data nodes to their the  $G_{/W}$  representative; the dr multimap stores the reverse mapping (summary node to all  $G$  node it represents). We call *class set* a set of types attached to a same resource in  $G$ ; to each non-empty class set corresponds exactly one node in  $G_{/TW}$  and  $G_{/TS}$ . clsd associates the corresponding summary node to each class set, while dc1s stores the class set of every (typed) data node from  $G$ . In a weak summary  $G_{/W}$ , every data property appears only once (Property 3); in  $G_{/TW}$ , appears *at most once with an untyped source and untyped target* (it may appear several times e.g., with subjects of different types and/or with typed vs. untyped subjects, see the four *title* edges in the  $G_{/TW}$  summary in Figure 15). The maps dc1s, dpSrc and dpTarg associate to a data property, its source and target (in  $G_W$ ), respectively, its (possible) single untyped source and untyped target (in  $G_{TW}$ ). Conversely, for each  $G$  data node (for  $G_W$ ), respectively, for each *untyped*  $G$  data node ( $G_{TW}$ ), the sets of data properties it is a source, respectively target of are stored in the maps srcDps (targDps). dtp maps every data property to the summary triple(s) it appears in. Each source and target clique is assigned an integer ID, and the IDs are stored in two lists, sc and tc. sToSc and oToTc associate to each data node, the ID of its source respectively target clique. Finally, for each source (target) clique ID, the  $G_{/S}$  (and  $G_{/TS}$ ) nodes having that clique are stored in scToSrc and tcToTarg.

Building  $G_{/W}$  and  $G_{/TW}$  requires a single pass over the data. We create candidate summary nodes as we go over the  $G$  triples, and record the association between  $G$  nodes and the respective summary nodes in the maps described above, gradually gaining knowledge about the data properties that each  $G$  node is a source and target of. When nodes share a source or target property, their summary representative nodes are fused (again based on Property 3). Building  $G_{/S}$  and  $G_{/TS}$  requires a first pass to compute the cliques, and a second to build a summary node for each (source clique, target clique) pair.

**Summarizing data triples in  $G_{/W}$ .** Algorithm 2 shows how we summarize data triples in a weak summary. The methods *getSource* and *getTarget* implement the representation functions, which map data nodes from  $G$  to data nodes in  $G_{/W}$ . These methods create the respective nodes the first time they are called for a

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

**Input** :  $s, p, G/W$   
**Output**: Data triples represented in  $G/W$   
Variables:  $src_u$  - data node representing an untyped source of  $p$ ;  
 $src_s$  - data node representing a (possibly typed) resource  $s$   
**Procedure**  $getSource(s, p, G/W)$

```

     $src_u \leftarrow dpSrc.get(p); src_s \leftarrow rd.get(s)$ 
    if  $src_u = \perp \wedge src_s = \perp$  then
    | return  $createDataNode(G/W, s)$ 
    else if  $src_u \neq \perp \wedge src_s = \perp$  then
    |  $rd.put(s, src), dr.put(src, s)$ 
    | return  $src_u$ 
    else if  $src_u = \perp \wedge src_s \neq \perp$  then
    | return  $src_s$ 
    else
    | if  $src_s = src_u$  then
    | | return  $src_s$ 
    | else
    | | return  $mergeDataNodes(G/W, src_s, src_u)$ 
    | end
    end
endProcedure
Procedure  $createDataNode(G/W, r)$ 
|  $d \leftarrow newInteger( ); rd.put(r, d); dr.put(d, r)$ 
| return  $d$ 
endProcedure

```

**Algorithm 3:** Representing the subject  $s$  of a data property  $p$  in  $G$  with a data node in  $G/W$

specific data property; on subsequent calls, the respective nodes are returned.

Algorithm 3 shows how we map the subjects of data triples in  $G$  to data nodes in  $G/W$ . After the method  $getSource$  has been executed, the source node of the property  $p$ , denoted  $src_u$ , and the node representing the subject  $s$ , denoted  $src_s$  must be the same. If neither  $src_u$  nor  $src_s$  exist yet,  $src$  will be a new data node representing  $s$  (line 3). In the cases when one of the nodes exist and the other does not, we simply use the existing node as  $src$  (lines 3-3). Moreover, if  $src_s$  had not existed, it means that  $s$  was unrepresented, so we assign  $src$  as its representative (line 3). On the other hand, if both  $src_u$  and  $src_s$  exist and they are the same, it does not matter which one we choose as  $src$  (line 3). When they differ, we have to merge them (line 3). 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 2 is very similar to  $getSource$ , the only difference being in passing the object  $o$  instead of the subject, and working with untyped property targets instead of sources.

Algorithm 4 outlines the summarization of type triples in  $G/W$ . It iterates over the subjects and classes of all type triples and tries to represent each resulting pair as follows. We look up the summary data node  $src$  representing  $s$ . If  $s$  is attached to any data property in  $G$ , this means  $s$  has already been represented when summarizing data triples and we assign its representative data node to  $src$  and we add the class to the class set of  $src$ . Otherwise,  $s$  has some types but no data property, thus we record it as a typed-only resource. Once all type triples whose subjects also had some data property have been represented, if there are any typed-only resources, we represent them as well. The procedure  $representTypedOnly$  creates the single data node  $N_\tau$  representing all resources from the list of typed-only resources and having *all* the types of typed-only resources (recall Section 7.1).

```

Input : Type triples table  $T_G$ , the summary  $G_W$ 
Output: Type triples represented in  $G_W$ 
 $toRes \leftarrow []$ ;  $toCls \leftarrow []$ 
foreach  $(s, c)$  in  $T_G$  do
  |  $repr \leftarrow representTypeTriple(G_W, s, c)$ 
  | if  $repr = false$  then
  | |  $toRes.add(s)$ ,  $toCls.add(c)$ 
  | end
end
if  $toRes.size > 0$  then
  |  $representTypedOnly(G_W, toRes, toCls)$ 
end
Procedure  $representTypeTriple(G_W, s, c)$ 
  |  $d \leftarrow rd.get(s)$ 
  | if  $d = \perp$  then
  | | return false
  | end
  |  $cls_d \leftarrow dcls.get(d)$ ,  $cls_d.add(d, c)$ 
  | return true
endProcedure
Procedure  $representTypedOnly(G_W, toRes, toCls)$ 
  |  $N_\tau \leftarrow newInteger()$ 
  | foreach  $r$  in  $toRes$  do
  | |  $rd.put(r, N_\tau)$ ,  $dr.put(N_\tau, r)$ 
  | end
  |  $cls_d \leftarrow dcls.get(N_\tau)$ 
  | foreach  $c$  in  $toCls$  do
  | |  $cls_d.add(c)$ 
  | end
endProcedure

```

**Algorithm 4:** Summarizing type triples in  $G_W$ 

### 9.3.1 Typed weak summary

**Summarizing data triples in  $G_{TW}$ .** The basic procedure for summarizing data triples in  $G_{TW}$  (Algorithm 5) is very similar to the one of weak presented in Algorithm 2, with the difference that the entries are stored to the maps only if  $src$  and  $target$  respectively, are untyped<sup>5</sup>. Further, the mapping of data nodes in  $G_{TW}$  is similar to the one of weak summary; they diverge when both  $src_s$  and  $src_u$  exist. *The decision on merging data nodes depends on the typing of  $src_s$ .* If  $src_s$  is typed, or it is untyped but already the same as  $src_u$ , we choose  $src_s$  as  $src$  (lines 5-5). On the other hand, if the two nodes are both untyped and different from each other, we must merge them (line 5). Therefore, in the typed weak summary only untyped data nodes may be merged, while the weak allows for the merging of typed data nodes as well. More precisely, since in the weak summary the data triples are represented *before* the type triples, the typing information is not yet available - the merged nodes may or may not be typed. However, in the typed weak summary the type triples are represented first and taken into account during the summarization of data triples to build a finer-grained summary.

**Summarizing type triples in  $G_{TW}$ .** Algorithm 6 shows how we summarize type triples in the typed weak summary. For each distinct subject we retrieve its class set (line 6) and the representative data node of

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

**Input** : The subject  $s$  of the data triple, the property  $p$  of  $s$ , the summary  $G_{/TW}$

**Output**: Data triples represented in  $G_{/TW}$

```

Procedure getSource( $s, p, G_{/TW}$ )
   $src_u \leftarrow dpSrc.get(p); src_s \leftarrow rd.get(s)$ 
  if  $src_u = \perp \wedge src_s = \perp$  then
    | return createDataNode( $G_{/TW}, s$ )
  else if  $src_u \neq \perp \wedge src_s = \perp$  then
    |  $rd.put(s, src), dr.put(src, s)$ 
    | return  $src_u$ 
  else if  $src_u = \perp \wedge src_s \neq \perp$  then
    | return  $src_s$ 
  else
    | if  $isTyped(src_s) \vee (src_s = src_u)$  then
      | | return  $src_s$ 
    | else
      | | return mergeDataNodes( $G_{/TW}, src_s, src_u$ )
    | end
  end

```

**endProcedure**

**Algorithm 5:** Representing the subject of a data property in  $G$  with a data node in  $G_{/TW}$

the class set (line 6). Observe that in a typed weak summary there is one data node *per distinct class set* which represents all subjects having the class set. If such a data node already exists for the given class set, we simply store it as the representative of the current subject (line 6). Otherwise, we create a new data node  $d$ , which represents  $s$  and has the same class set as  $s$  (line 6). Further,  $d$  is stored as the representative of its class set (line 6).

**Input** : Type triples table  $T_G$ , the summary  $G_{/TW}$

**Output**: Type triples represented in  $G_{/TW}$

$typ \leftarrow eval(SELECT s, c FROM T_G ORDER BY s)$

```

foreach distinct  $s$  in  $typ$  do
  |  $cls_s \leftarrow dcls.get(s)$ 
  |  $d \leftarrow clsd.get(cls_s)$ 
  | if  $d \neq \perp$  then
    | |  $rd.put(s, d), dr.put(d, s)$ 
  | else
    | |  $d \leftarrow createDataNodeG_{/TW}, cls_s, s$ 
    | |  $clsd.put(cls_s, d)$ 
  | end
end

```

**end**

**Procedure** createDataNode( $G_{/TW}, r, cls_r$ )

```

  |  $d \leftarrow newInteger, rd.put(r, d), dr.put(d, r), dcls.put(d, cls_r)$ 
  | return  $d$ 

```

**endProcedure**

**Algorithm 6:** Summarizing type triples in  $G_{/TW}$

**Summarizing schema.** Prior to summarization, the encoded schema table was created. Since the schema of  $G$  and  $G_{/TW}$  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.

### 9.3.2 Strong summary

**Summarizing data triples in  $G_{/s}$ .** Instead of building the cliques gradually as in the weak summary, in Algorithm 7 we compute all source and target cliques at the very beginning (lines 1-2). Then in lines 3-15, for each data node  $d$  in  $G$ , a representative data node  $d_{repr}$  is assigned in  $G_{/s}$ . This data node  $d_{repr}$  must have the exact same source and target clique as  $d$ , and there can be at most one such node for any two source and target cliques. Finally, for each data triple  $s p o$  in  $D_G$ , we retrieve the representative data node of  $s$ , denoted  $src$ , and of  $o$ , denoted  $targ$ , forming a data triple  $src p targ$  in  $G_{/s}$ . The method for creating the data triple is the same as for the weak summary in Algorithm 2.

**Input** : Data triples table  $D_G$ , the summary  $G_{/s}$

**Output**: Data triples represented in  $G_{/s}$

Variables:  $scId, tcId$  - source (target) clique ID;  $D_{SC}, D_{TC}$  - the set of source (target) data nodes in  $G_{/s}$  representing the source clique  $scId$  (target clique  $tcId$ )

buildSourceCliques( $D_G, G_{/s}$ )

buildTargetCliques( $D_G, G_{/s}$ )

**foreach** *distinct data node*  $d$  in  $D_G$  **do**

$scId \leftarrow sToSc.get(d), tcId \leftarrow oToTc.get(d)$

$D_{SC} \leftarrow scToSrc.get(scId), D_{TC} \leftarrow tcToTarg.get(tcId)$

**if**  $D_{SC} \cap D_{TC} \neq \emptyset$  **then**

$rd.put(r, d_{repr}), dr.put(d_{repr}, r)$

**else**

$d_{repr} \leftarrow createDataNode(G_{/s}, d)$

$scToSrc.put(scId, d_{repr}), tcToTarg.put(tcId, d_{repr})$

**end**

**end**

**foreach**  $s p o$  in  $D_G$  **do**

$src \leftarrow rd.get(s), targ \leftarrow rd.get(o)$

**if**  $!existsDataTriple(G_{/s}, src, p, targ)$  **then**

$createDataTriple(G_{/s}, src, p, targ)$

**end**

**end**

**Algorithm 7:** Summarizing data triples in  $G_{/s}$

The procedure buildSourceCliques( $D_G, G_{/s}$ ) invokes the procedure computeSourceClique( $G_{/s}, s, P_s$ ) shown in Algorithm 8 for each subject in  $D_G$  and the set of data properties of which it is the source.

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

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

**Summarizing type triples in  $G_{/s}$ .** The types in  $G_{/s}$  are summarized in the same manner as for the  $W_G$ , described in Algorithm 4.

**Property 8.** (ALGORITHM CORRECTNESS) *Given an RDF graph  $G$ : Algorithms 2 and 4, together, build the weak summary  $G_{/w}$ . Algorithms 5 and 6 build the typed weak summary  $G_{/TW}$ . Algorithm 7 shows how to build the strong summary  $S_G$ ; building  $TS_G$  is very similar.*



**Input** : The summary  $G_{/s}$ , subject  $s$  from  $D_G$ , the set of data properties  $P_s$  of which  $s$  is the source

**Output**:  $sc$  list and  $sToSc$  map of  $G_{/s}$  populated

**Procedure** computeSourceClique( $G_{/s}$ ,  $s$ ,  $P_s$ )

```

    srcClq ← ⊥
    foreach c in sc do
        if c ∩ P_s ≠ ∅ then
            srcClq ← c
            break;
        end
    end
    if srcClq ≠ ⊥ then
        srcClq.addAll(P_s)
    else
        srcClq.addAll(P_s)
        sc.add(srcClq)
    end
    sToSc.put(s, sc.indexOfsrcClq+1)
endProcedure

```

**Algorithm 8:** Computing source cliques in  $G_{/s}$

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

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

## 10 Summary comparison

How do our summaries relate to each other, and how can one compare them?  $G_{/w}$ ,  $G_{/s}$ ,  $G_{fw}$ ,  $G_{bw}$  and  $G_{fb}$  are complete (Theorems 2, 3 and 4), while  $G_{/TW}$  and  $G_{/TS}$  are not (Propositions 5 and 7).

**Composing summaries.** One can easily prove that  $(G_{/s})_{/w} = G_{/w}$ , i.e., one could compute  $G_{/w}$  by first summarizing  $G$  into  $G_{/s}$ , and then applying weak summarization on this (typically much smaller) graph; similarly, one can show that  $(G_{/TS})_{/TW} = G_{/TW}$ . It is also the case that  $(G_{/w})_{/s} = G_{/w}$ , i.e., strong summarization cannot compress a weak summary further, and similarly  $(G_{/TW})_{/TS} = G_{/TW}$ . Figure 20 summarizes the main relationships between summaries.

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

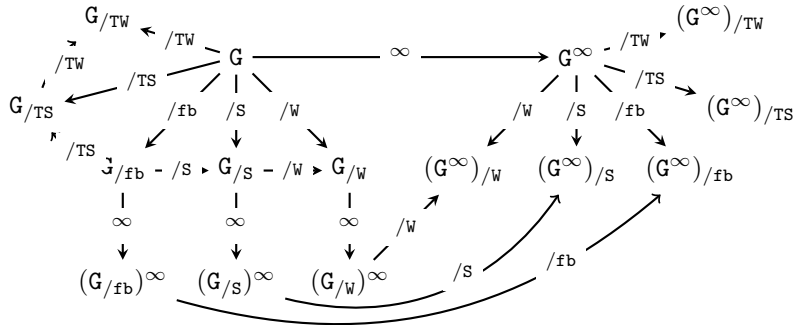


Figure 20: Main relations between our summaries.

summarizes at least as much as the typed strong summary  $G_{/TS}$ .

**Accuracy comparison.** Another interesting measure is the closeness between a graph and its summary. By Definition 7, any subgraph of the  $G$  can be embedded through homomorphism in any summary of  $G$ ; however, the opposite is not true, as the representation of data nodes by summary nodes may lead to subgraphs in the summary which do not exist in  $G$ . For instance,  $G_{/W}$  in Figure 13 a node has outgoing edges labeled  $a$  and  $e$ , whereas such a node does not exist in the original graph (Figure 10). Formally:

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

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

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

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

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

**Relation between clique- and bisimulation-based summaries.** Our  $W$  and  $S$  summaries are symmetric w.r.t.  $G$  edge source and target nodes, thus the comparison with the unidirectional  $FW$  and  $BW$  bisimulation summaries is limited. A deeper comparison is possible and interesting with the (symetric)  $FB$ -bisimulation summary  $G_{/fb}$ . Since  $G_{/fb}$  preserves  $G$  structure and thus node cliques, we have  $(G_{/fb})_{/W} = G_{/W}$  and  $(G_{/W})_{/fb} = G_{/W}$ , and similar statements for  $G_{/S}$ ,  $G_{/TW}$ , and  $G_{/TS}$ . Different from our the clique-based,  $G_{/fb}$  guarantees  $AL = 0$ , but it is usually *much* larger, as illustrated in Section 6.2 and in our experiments, making it oftentimes impractical.

## 11 Experiments

We report here on experiments based on our algorithms described in Section 9.

G	G	G <sup>∞</sup>	G <sub>/≡<sub>fw</sub></sub>	G <sup>∞</sup> <sub>/≡<sub>fw</sub></sub>	S <sub>G</sub>	t (s)	t <sub>S</sub> (s)	Speed-up
BSBM 1M	1,000,708	1,009,138	7,124	8,076	150	47.95	44.71	<b>6.75%</b>
BSBM 10M	10,538,484	10,628,484	16,883	18,833	584	487.54	446.87	<b>8.34%</b>
LUBM 1M	1,227,868	1,388,296	1,989	1,989	124	69.81	62.77	<b>10%</b>
LUBM 10M	11,990,183	13,562,978	25,002	25,002	124	679.29	599.39	<b>11.76%</b>
INSEE Geo	369,042	1,126,239	12,524	26,954	281	145.14	67.59	<b>53.43%</b>

Table 2: Summarization experiments.

**Experiments settings.** We stored  $G$  in PostgreSQL v9.6 (encoding URIs and literals as integers for efficiency), while the bisimulation summary is stored in Apache MapDB ([www.mapdb.org](http://www.mapdb.org)). We used a server with an Intel Xeon CPU E5-2640 v4 @2.40GHz and 124 GB of memory; Postgres had 30 GB of shared buffers and 640 MB working memory. The JVM had 70 GB of RAM.

**Datasets.** We used *synthetic* RDF graphs (BSBM [2] and LUBM [13] benchmark data) as well as *real-life* graphs: INSEE Geo geographic data (<http://rdf.insee.fr/geo/index.html>) and DBLP. The numbers of data triples, respectively type triples in  $G^\infty$  (denoted  $G_{\text{sat}}$  data and  $G_{\text{sat}}$  type, respectively) appear at the top of Figure 21; note the logarithmic vertical axis. We do not show the  $G^\infty$  schema size as it is much smaller (in the hundreds) in all cases.

**Baseline: bisimulation summaries.** To compare our work, we picked *bisimulation-based RDF summaries* to instantiate in our RDF summarization framework. We chose these as a comparison baseline because (i) they are *the only representative summaries (Definition 1) taking into account the semantics and types of RDF, that we are aware of, beyond the four summaries introduced here*; (ii) bisimulation summaries have been extensively studied in the literature. We omit  $G_{/b_w}$  results for space limitations and because they are between those of  $G_{/f_b}$  and  $G_{/f_w}$  (see below).

**Bisimulation summarization and shortcut experiments.** We built  $G_{f_b}$ ,  $G_{f_w}$  and  $G_{b_w}$  for synthetic [2, 13] and real<sup>6</sup> RDF graphs, outlined in Table 2;  $|\cdot|$  denotes the size (number of triples) of an RDF graph. We found that  $G_{/f_b}$  can be almost as large as  $G$ , e.g., for BSBM 1M, summarization reduces the number of triples by only 0.02%! As  $G_{/f_w}$  is often much smaller than  $G$ , we only report on experiments with  $G_{/f_w}$ . In Table 2,  $t$  is the time to saturate then summarize  $G$ , while  $t_S$  is the time to compute the same result following the shortcut (green path in Figure 7). The shortcut is faster in all cases, and significantly so (53.43%) for the INSEE dataset where  $G^\infty$  is much larger than  $G$ ; recall also that the shortcut is the only option when one cannot saturate  $G$ . Further, we see that  $(G^\infty)_{/f_w}$  is 41 to 697 times smaller than  $G^\infty$ , which confirms that the summary is much easier to handle than  $G$ ; interactive graph visualization techniques and/or querying (e.g., in the example from the Introduction, ask “find the politicians owning companies”) allow to determine the usefulness of  $G$  based on its representative summary  $(G^\infty)_{/f_w}$ .

**Summaries size.** Figure 21 (bottom) shows the number of data, types, and schema triples for  $(G^\infty)_{/w}$  and  $(G^\infty)_{/s}$  (the size of the saturated schema is the same). We define the **compression factor**  $cf_{\equiv}$  as  $|G|/|G_{\equiv}|$ , the ratio between the numbers of triples in  $G$  and in the summary. In our experiments, the  $G_{/T_w}$  is very close or equal to the size of  $G_{/w}$ , while the  $G_{/T_s}$  is close to that of  $G_{/s}$ . Compression factors for the same graphs as above, for our weak and strong summaries as well as for  $G_{f_b}$  and  $G_{f_w}$ , also appear in Figure 22. Figure 22 shows that the compression factor for our four summaries is at least 500, often above  $10^3$ , and up to  $2 \cdot 10^6$  for DBLP dataset; in our experiments,  $cf_w$  was typically the highest. In contrast, since  $\equiv_{f_b}$  is a very strong condition,  $G_{/f_b}$  summarizes very little ( $cf_{f_b}$  rounds to 1 up to 3). Since  $\equiv_{f_w}$  is a weaker condition,  $G_{/f_w}$  compresses more ( $cf_{f_w}$  between 29 and 617), but still much less than our four summaries.

A more detailed look at BSBM summarization is provided in Figure 23, which shows node and edge

<sup>6</sup>INSEE Geographic dataset, available at: <http://rdf.insee.fr/geo/index.html>

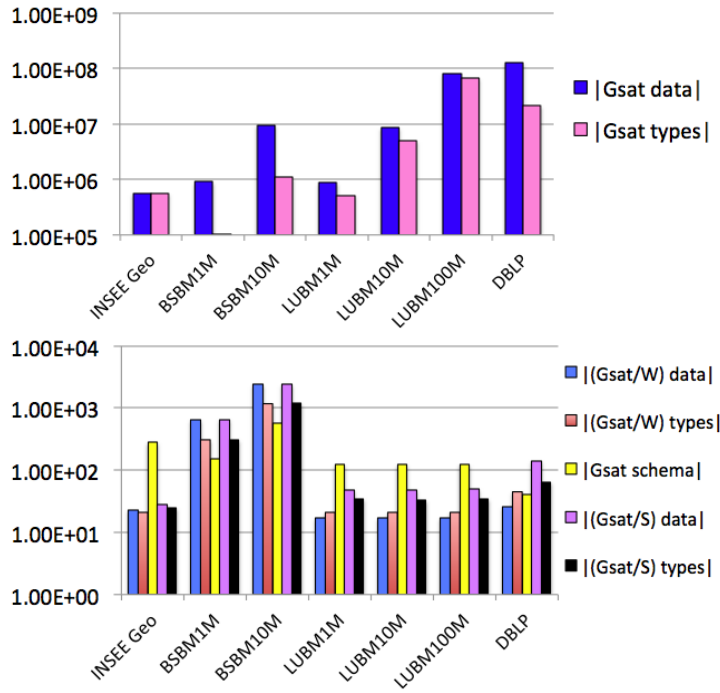


Figure 21: Top: data and type triple counts in  $G^\infty$  for various RDF graphs; bottom: data, type and schema triple counts in  $(G^\infty)_w$  and  $(G^\infty)_s$  for these graphs.

Graph $G$ and $G_{\equiv}$	$cf_{\equiv}$	$t$ (s)	$t_S$ (s)	Speed-up
INSEE Geo, $G_w$	1,186	21.42	0.66	<b>96.92%</b>
INSEE Geo, $G_s$	1,164	22.60	1.70	<b>92.48%</b>
INSEE Geo, $G_{fw}$	29	148.83	67.76	<b>54.47%</b>
INSEE Geo, $G_{fb}$	3	310.54	270.03	<b>13.04%</b>
BSBM1M, $G_w$	956	3.02	1.30	<b>56.95%</b>
BSBM1M, $G_s$	937	5.51	3.53	<b>35.94%</b>
BSBM1M, $G_{fw}$	140	61.68	58.71	<b>5.81%</b>
BSBM1M, $G_{fb}$	1	332.95	628.23	<b>-88.68%</b>
LUBM1M, $G_w$	7579	7.56	1.39	<b>81.61%</b>
LUBM1M, $G_s$	5903	9.49	3.11	<b>67.23%</b>
LUBM1M, $G_{fw}$	617	86.21	78.74	<b>8.66%</b>
LUBM1M, $G_{fb}$	1	314.82	621.14	<b>-97.29%</b>
LUBM100M, $G_w$	705,897	1633.94	144.65	<b>91.15%</b>
LUBM100M, $G_s$	547,155	1810.26	557.95	<b>69.18%</b>
DBLP, $G_w$	2,123,767	1180.04	238.96	<b>79.75%</b>
DBLP, $G_s$	731,978	1325.06	630.89	<b>52.39%</b>
DBLP, $G_{fw}$	574	20521.23	19363.34	<b>5.64%</b>

Figure 22: Compression factors, time to generate  $(G^\infty)_w$  and  $(G^\infty)_s$  directly and through the shortcut for a variety of RDF graphs, and the shortcut speed-up.

counts for summaries of BSBM graphs of various sizes. The horizontal axis is labeled in input triples; the vertical axis is in log scale. We plot the number of class nodes of  $G$  next to the summary data nodes, and the number of schema edges next to the edge counts. The figure shows that  $G_w$  and  $G_s$  strongly summarize graph structure, while  $G_{fw}$  and  $G_{fb}$  which isolate typed from untyped data nodes lead to larger, more complex summaries.

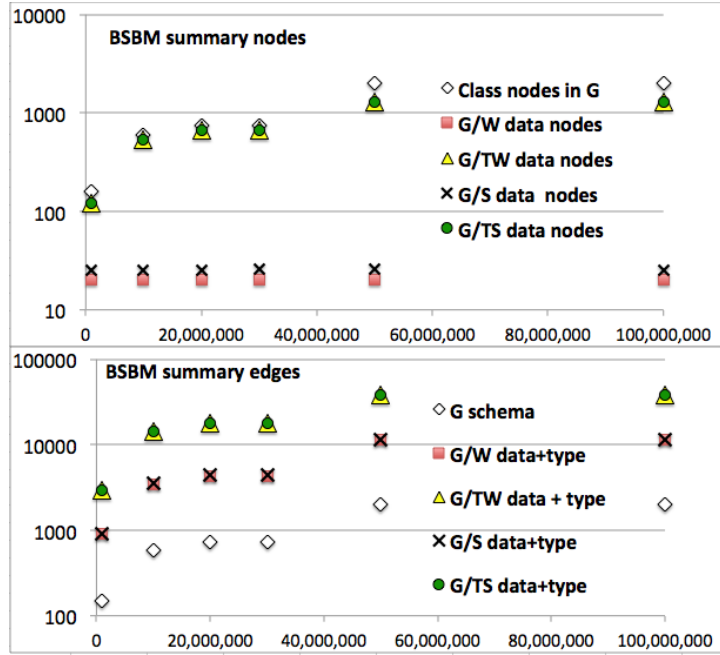


Figure 23: Summary node counts (top) and edge counts (bottom) for BSBM graphs.

**Clique summarization time and scale-up.** Figure 24 shows the average time (in seconds, 3 executions) to build clique summaries; they go up to 800 seconds for 100M triples. Times for the LUBM and BSBM synthetic data are also plotted in the lower part of the figure: both axes are in logarithmic scale, the horizontal axis labeled in  $|G|$  triples). The graphs show that summarization time grows almost linearly, and the time to build  $G_{/W}$  grows fastest. This is because building  $G_{/W}$  and  $G_{/S}$  maximizes the number of nodes for which cliques are analyzed and fused; in contrast, when building  $G_{/TW}$  and  $G_{/TS}$ , *the cliques of typed data nodes are never constructed (directly or indirectly)*.  $G_{/S}$  and  $G_{/TS}$  often take longer to build than  $G_{/W}$  and  $G_{/TW}$ ; this is because of the extra effort needed to build the cliques, while  $G_{/W}$  and  $G_{/TW}$  use a more direct, faster method. We find summarization times acceptable, also considering that this is an off-line task; if speed-up is important, our algorithms can easily be parallelized to take advantage of a parallel processor (as in [19]) or a parallel computing framework such as Spark.

**Shortcut experiments.** Figure 22 also shows the time to build  $(G^\infty)_{/\equiv}$  in two ways: (i) saturate  $G$  then summarize, denoted  $t$ , and (ii) summarize  $G$ , saturate the summary, and summarize again, denoted  $t_S$ . We define the shortcut **speed-up** as  $(t - t_S)/t$  and report it also in the figure. The speed-up is highest for  $G_{/W}$  (up to almost 97%) and  $G_{/S}$  (up to 97%): this is a direct consequence of their good compression. Indeed,  $t$  includes the time to summarize  $G^\infty$ , while  $t_S$  includes the time to summarize  $(G_{/\equiv})^\infty$ . If the latter is much smaller,  $t_S$  is smaller and the speed-up is positive. The much more modest compression factor of  $G_{/\equiv W}$  translates into a modest speed-up (5% to 13%), while for  $G_{/\equiv b}$ , the speed-up is negative! This is because with a basically unnoticeable compression through  $\equiv_{fb}$ , there is no time gain to summarize  $(G_{/\equiv b})^\infty$  instead of  $G^\infty$ . Thus, the overhead of the *three* shortcut steps (instead of two on the other path) is not offset, that is,  $t < t_S$ .

**Experiment conclusion.** Our experiments have shown that our four summaries can be built efficiently, and they reduce graph sizes by factors hundreds up to  $2 \cdot 10^6$ ; they are *two to three orders of magnitude more compact than bisimulation-based summaries*;  $G_{/W}$  was the most compact and representative RDF summary, and the other three summaries we introduced are close. Finally, among the summaries which admit a shortcut ( $G_{/W}$ ,  $G_{/S}$ ,  $G_{/\equiv b}$  and  $G_{/\equiv W}$ ), the shortcut speed-up is up to 96% for  $G_{/W}$ , trailed by  $G_{/S}$ , while

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

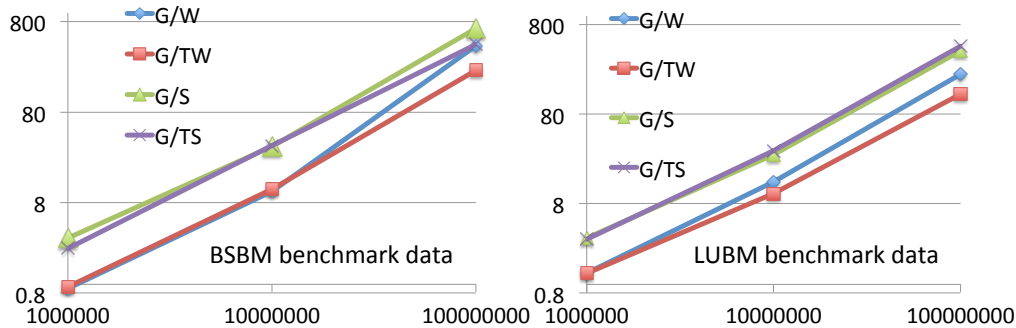


Figure 24: Summarization times (s) and their scale-up with the input size, for various RDF datasets.

for  $G_{fw}$  the speed-up is small and for  $G_{fb}$  it is negative (the shortcut has no interest here). Despite their polynomial complexity, in practice, our algorithms scale basically linearly. Overall, if summary size or building time are a critical concern, we recommend using  $G_w$ ; otherwise,  $G_s$  may give finer-granularity information.

## 12 Related Work

Graph summarization has a long history and many applications.

**Non-RDF graph summaries** ignore graph semantics, thus representativeness, a main benefit of our summaries, is not guaranteed. Dataguides [12] (including “strong” ones) are not quotients, as a graph node may be represented by *more than one* dataguide node, thus a dataguide may be larger than the original graph, and its construction has worst-case exponential time complexity in the size of  $G$ . [24, 8, 16, 10, 23] build quotient summaries based on (sometimes unidirectional, sometimes  $k$ -bounded) bisimulation, which, as shown in Section 4.1, prevents representativeness. Further, we observed and confirmed by our experiments in the preceding section, that such summaries do not achieved significant graph compression. [11] considers reachability and graph pattern queries on labeled graphs, and builds *answer-preserving summaries* for such queries; this property subsumes representativeness. However, their query semantics is not based on graph homomorphisms, which underlie SPARQL and the dialects we use, but on *bounded graph simulation*. Under this semantics, answering a query becomes  $P$  (instead of  $NP$ ), at the price of not preserving the query structure (i.e., joins), which our representativeness preserves (Proposition 1).

More distant from our work, [7] introduces an *aggregation* framework for OLAP on labeled graphs,

whereas we focus on representing graph structure and semantics. [6] builds a set of randomized summaries to be *mined* instead of the original graph for better performance, with guaranteed bounds on the information loss. Graph compression with bounded “error” (number of edges to be added as “corrections” after decompression, to retrieve the original graph) is studied in [25, 17]. Quantitative summarization, where nodes and edges are summarized by their counts, is the focus of [21]. [22] surveys many other quantitative, mining-oriented graph sampling and summarization methods.

**RDF summarization** [9, 29, 30] study bisimulation-based RDF quotient summaries, in particular providing efficient parallel summarization algorithms [9] and showing they are representative for data-only queries [29]. However, these summaries ignore the special role type and schema triples play in RDF, and therefore cannot guarantee representativeness for the implicit data and type-aware queries. We have included bisimulation-based RDF summaries in our framework, and showed they admit our shortcut.

[3] introduces a family of summaries for RDF graphs including data and type triples (not schema triples). The authors define a set of quality metrics quantifying a summary’s loss of precision, and compare their efficiency, effectiveness and precision based on a MapReduce implementation of the summarization algorithm. These summaries are also based on node equivalence relations, including forward-, backward- and forward-and-backward bisimilarity which we used to illustrate our framework here.

Importantly, *the summary of [3] is not a quotient (Definition 4) whereas our summary is a quotient through an RDF equivalence relation (Definition 7)*. Indeed, in [3] there may exist graph nodes “for which there is no  $\sim$ -equivalent class” (Definition 2.4); the authors capture them all into *a single equivalence class* named BLANK  $\sim$ -EQUIVALENCE CLASS. In contrast, in our framework, any node  $n$  which is not equivalent to any other node:

- If  $n$  is a class or property node, Definition 6 ensures it is represented by itself; this is crucial for our summary to preserve the input schema (Property 1), and for our shortcut to hold;
- If  $n$  is a data node, its equivalence class (of size 1) is represented in the summary by a node with a *fresh* URI (Definition 7), in particular different from all the other summary nodes.

Thus, we never put nodes not equivalent to any other in the same equivalence class, as in [3]. Even restricted to just data triples, our summary and the one of [3] may lead to different results.

Also because the summary of [3] is not a quotient, the summarization they consider may assign to “several equivalence classes to a node” (Section 4.1, second paragraph in the column at right, in particular for their relation  $\sim_{st}$ ). In contrast, since our summary (Definition 7) is a quotient, each node is always assigned to exactly one equivalence class and represented by exactly one node.

Another core difference between our framework and the one of [3] is that our shortcut *summarizes both the data and the (possible) RDF Schema of  $G$* , that is, the *semantics* of  $G$ . Also, the shortcut result is *the same for any graph  $G$  having the saturation  $G^\infty$* . In contrast, the summaries introduced in [3] which use the type triples would have to be recomputed if the set of type triples change. Or, type triples are naturally added when one adds entailed triples to  $G$ .

[27] introduces a *simulation* RDF quotient based on *triple (not node) equivalence*. (Intuitively, simulation is a unidirectional “half” of bisimulation, e.g., replace “iff” by “if” in the bisimilarity condition). For compactness, they bound the simulation; this (and other aspects of their proposal) prevents representativeness.

Other proposals are based on clustering [14] user-defined aggregation rules [28], mining [6], and identification of frequent subtrees [35]; they are not representative in our sense either. Finally, some summaries record structural or statistical aspects about an RDF graph: [18] uses forward bisimulation to characterize usage frequency of classes and properties from some namespaces, while [26] computes statistics reflecting data property frequency, the most specific types of properties’ subjects/objects etc. None of the abovementioned RDF summary proposals ensures representativeness.

We had demonstrated [5] and (informally) presented  $G_{/W}$  and  $G_{/TW}$  in [4], with procedural definitions (not as quotients); unfortunately, the shortcut described there was incorrect. This report formalizes our RDF summarization framework, shows how to instantiate it for the bisimulation-based summaries  $G_{fw}$ ,  $G_{bw}$  and  $G_{fb}$ , presents the first formalization of  $G_{/W}$  and  $G_{/TW}$ , and introduces the novel  $G_{/S}$  and  $G_{/TS}$ . Further, we show that  $G_{fw}$ ,  $G_{bw}$  and  $G_{fb}$ ,  $G_{/W}$  and  $G_{/S}$  admit our shortcut, and experimentally demonstrate that the latter two are several orders of magnitude more compact than the (comparable) summaries  $G_{/fw}$  and  $G_{/fb}$ ; the shortcut also translates into dramatic speed-up factors for summarizing  $G^\infty$ . This demonstrates their strong theoretical properties and practical interest.

## 13 Conclusion

We studied the problem of building summaries of RDF graphs, accounting for their structure and semantics; in particular, we focused on *quotient* graph summaries, widely studied in the literature. We defined a framework for representative summarization of RDF graph, based on RDF equivalence relations carefully crafted to reflect the particular role schema and types play in RDF. The complete semantics of an RDF graph under the RDF Schema entailment regime is its saturation, thus the most meaningful summary for a given graph  $G$  is that of  $G^\infty$ . To obtain  $(G^\infty)_{/\equiv}$  even without saturating  $G$ , we exhibit a shortcut (Theorem 1) and provide a sufficient condition for the shortcut to hold on a given RDF summary. Further, we introduced a new concept of *property cliques* out of which we build four novel *clique summaries*, two of which (as well as bisimulation-based ones) can be *efficiently* computed through our shortcut method. All our clique summaries are *representative* w.r.t. a large and useful dialect of SPARQL; they can be built efficiently and achieve RDF graph compression factors of hundreds to  $2 \cdot 10^6$ ; the highest values are reached for real-life datasets. Future work will improve scalability by leveraging a massively parallel platform such as Spark. We are also interested in devising advanced visualizations of our RDF summaries.

## References

- [1] RDF Summary project website: <https://team.inria.fr/cedar/projects/rdfsummary>, 2017.
- [2] 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] Š. Čebirić, F. Goasdoué, and I. Manolescu. Query-oriented summarization of RDF graphs. In *BICOD*, 2015.
- [5] Š. Čebirić, F. Goasdoué, and I. Manolescu. Query-oriented summarization of RDF graphs (demonstration). *PVLDB*, 8(12), 2015.
- [6] C. Chen, C. X. Lin, M. Fredrikson, M. Christodorescu, X. Yan, and J. Han. Mining graph patterns efficiently via randomized summaries. *PVLDB*, 2(1), 2009.
- [7] C. Chen, X. Yan, F. Zhu, J. Han, and P. S. Yu. Graph OLAP: towards online analytical processing on graphs. In *ICDM*, 2008.
- [8] Q. Chen, A. Lim, and K. W. Ong.  $D(K)$ -index: An adaptive structural summary for graph-structured data. In *SIGMOD*, 2003.
- [9] M. P. Consens, V. Fionda, S. Khatchadourian, and G. Pirrò. S+EPPs: Construct and explore bisimulation summaries + optimize navigational queries; all on existing SPARQL systems (demonstration). *PVLDB*, 8(12), 2015.
- [10] M. P. Consens, R. J. Miller, F. Rizzolo, and A. A. Vaisman. Exploring XML web collections with DescribeX. *ACM TWeb*, 4(3), 2010.
- [11] W. Fan, J. Li, X. Wang, and Y. Wu. Query preserving graph compression. In *SIGMOD*, 2012.



- [12] R. Goldman and J. Widom. Dataguides: Enabling query formulation and optimization in semistructured databases. In *VLDB*, 1997.
- [13] Y. Guo, Z. Pan, and J. Hefflin. LUBM: A benchmark for OWL knowledge base systems. *J. Web Sem.*, 3(2-3), 2005.
- [14] S. Gurajada, S. Seufert, I. Miliaraki, and M. Theobald. Using graph summarization for join-ahead pruning in a distributed RDF engine. In *SWIM*, 2014.
- [15] M. R. Henzinger, T. A. Henzinger, and P. W. Kopke. Computing simulations on finite and infinite graphs. In *FOCS*, 1995.
- [16] R. Kaushik, P. Bohannon, J. F. Naughton, and H. F. Korth. Covering indexes for branching path queries. In *SIGMOD*, 2002.
- [17] K. Khan, W. Nawaz, and Y. Lee. Set-based approximate approach for lossless graph summarization. *Computing*, 97(12):1185–1207, 2015.
- [18] S. Khatchadourian and M. P. Consens. ExpLOD: Summary-based exploration of interlinking and RDF usage in the linked open data cloud. In *ESWC*, 2010.
- [19] S. Khatchadourian and M. P. Consens. Constructing bisimulation summaries on a multi-core graph processing framework. In *GRADES*, 2015.
- [20] W. Le, F. Li, A. Kementsietsidis, and S. Duan. Scalable keyword search on large RDF data. *IEEE TKDE*, 26(11), 2014.
- [21] K. LeFevre and E. Terzi. GraSS: Graph structure summarization. In *SDM*, 2010.
- [22] S.-D. Lin, M.-Y. Yeh, and C.-T. Li. Sampling and summarization for social networks (tutorial), 2013.
- [23] Y. Luo, G. H. L. Fletcher, J. Hidders, Y. Wu, and P. D. Bra. External memory k-bisimulation reduction of big graphs. In *CIKM*, 2013.
- [24] T. Milo and D. Suciu. Index structures for path expressions. In *ICDT*, 1999.
- [25] S. Navlakha, R. Rastogi, and N. Shrivastava. Graph summarization with bounded error. In *SIGMOD*, 2008.
- [26] M. Palmonari, A. Rula, R. Porrini, A. Maurino, B. Spahiu, and V. Ferme. ABSTAT: linked data summaries with abstraction and statistics. In *ESWC Workshops*, 2015.
- [27] F. Picalausa, Y. Luo, G. H. L. Fletcher, J. Hidders, and S. Vansummeren. A structural approach to indexing triples. In *ESWC*, 2012.
- [28] M. Rudolf, M. Paradies, C. Bornhövd, and W. Lehner. SynopSys: large graph analytics in the SAP HANA database through summarization. In *GRADES*, 2013.
- [29] A. Schätzle, A. Neu, G. Lausen, and M. Przyjaciół-Zablocki. Large-scale bisimulation of RDF graphs. In *SWIM*, 2013.
- [30] T. Tran, G. Ladwig, and S. Rudolph. Managing structured and semistructured RDF data using structure indexes. *IEEE TKDE*, 25(9), 2013.
- [31] T. Tran, H. Wang, S. Rudolph, and P. Cimiano. Top-k exploration of query candidates for efficient keyword search on graph-shaped (RDF) data. In *ICDE*, 2009.
- [32] W3C. Resource description framework. <http://www.w3.org/RDF/>.
- [33] W3C. SPARQL 1.1 query language. <http://www.w3.org/TR/sparql11-query/>, March 2013.
- [34] <http://gromgull.net/blog/2010/09/btc2010-basic-stats>, 2010.
- [35] P. Zhao, J. X. Yu, and P. S. Yu. Graph indexing: Tree + delta  $\geq$  graph. In *VLDB*, 2007.



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

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

Publisher  
Inria  
Domaine de Voluceau - Rocquencourt  
BP 105 - 78153 Le Chesnay Cedex  
[inria.fr](http://inria.fr)

ISSN 0249-6399