



Compact Summaries of Rich Heterogeneous Graphs

Šejla Čebirić, François Goasdoué, Paweł Guzewicz, Ioana Manolescu

**RESEARCH
REPORT**

N° 8920

June 2016

Project-Teams CEDAR



Compact Summaries of Rich Heterogeneous Graphs

Šejla Čebirić, François Goasdoué, Paweł Guzewicz, Ioana Manolescu

Project-Teams CEDAR

Research Report n° 8920 — version 6 — initial version June 2016 — revised
version July 2018 — 40 pages

Abstract:

Large data graphs with complex and heterogeneous structure, possibly featuring typed data and an ontology encoding the application-domain semantics, are widely used nowadays. The literature provides many solutions for building succinct representations of graphs, called summaries, in particular based on graph quotients through an equivalence relation between graph nodes.

We consider efficient and compact summarization of rich heterogeneous graphs, in particular RDF ones, which may feature *data edges*, *typed nodes*, and an *ontology*. First, we devise *new graph node equivalence relations*, particularly tolerant of structural heterogeneity; they lead to compact yet informative quotient summaries. Second, we show how to *extend any node equivalence relation* (including, but not limited to ours) to types and ontologies, and provide the first in-depth study of the *interplay between quotient summarization and RDF graph saturation*, which defines the semantics of an RDF graph, in particular in the presence of an ontology. We establish a sufficient condition on a node equivalence relation, which if met allows an efficient method, called *shortcut*, for summarizing RDF graphs. We describe novel, efficient, incremental algorithms for summarizing graphs with our node equivalence relations, and experiments validating their performance.

Key-words: Reasoning and Knowledge Representation, Databases, RDF Graphs, Semantic Web, Graph Summaries

RESEARCH CENTRE
SOPHIA ANTIPOLIS – MÉDITERRANÉE

2004 route des Lucioles - BP 93
06902 Sophia Antipolis Cedex

Compact Summaries of Rich Heterogeneous Graphs

Résumé : Les grands graphes de données, à structure potentiellement complexe et hétérogène, comprenant parfois des noeuds typés et/ou une ontologie décrivant la sémantique d'un domaine d'application, sont utilisés largement de nos jours. Des nombreuses méthodes pour résumer de tels graphes ont été proposées; chacune aboutit à une structure compacte représentant l'essentiel de l'information sur la structure du graphe. Une famille de telles méthodes est basée sur la construction d'un graphe quotient, à partir d'une relation d'équivalence entre les noeuds du graphe d'origine.

Dans ce travail, nous nous intéressons à la construction efficace de résumés compacts pour des graphes de données hétérogènes, en particulier des graphes RDF, qui peuvent comporter des données, des noeuds typés, et une ontologie. Nous proposons de nouvelles relations d'équivalence entre les noeuds d'un graphes, relations particulièrement adaptées à la hétérogénéité souvent rencontrée dans des graphes RDF. Ces relations conduisent à des résumés quotients qui sont compacts, tout en préservant de l'information sur la structure du graphe. Nous montrons comment tout relation d'équivalence entre les noeuds d'un graphe (comprenant, mais ne se limitant pas aux relations que nous introduisons) peut être étendue aux noeuds comportant des types, ainsi qu'aux ontologies. Par la suite, nous présentons la première étude approfondie de l'interaction entre les résumés par quotient et la saturation d'un graphe RDF, qui définit sa sémantique en présence d'une ontologie. Nous identifions une condition suffisante sur une relation d'équivalence entre des noeuds du graphe, condition qui, lorsqu'elle est satisfaite, conduit à une méthode rapide, appelée *shortcut* (raccourci) pour la construction du résumé d'un graphe RDF. Nous présentons des nouveaux algorithmes efficace, en particulier des algorithmes incrémentaux, pour résumer des graphes par nos nouvelles relations d'équivalence. Enfin, nous présentons des expériences qui valide la performance de nos algorithmes et l'intérêt des résumés étudiés.

Mots-clés : Représentation de connaissances, Bases de Données, RDF, Web sémantique, Résumés de Graphes

Contents

| | | |
|-----------|---|-----------|
| 1 | Introduction | 4 |
| 2 | Data graphs | 5 |
| 3 | Summarization framework | 7 |
| 4 | Data graph summarization | 9 |
| 4.1 | Data property cliques | 9 |
| 4.2 | Strong and weak node equivalences | 10 |
| 4.3 | Weak and strong summarization | 11 |
| 5 | Typed data graph summarization | 12 |
| 5.1 | Data-then-type summarization | 13 |
| 5.2 | Type-then-data summarization | 13 |
| 6 | RDF graph summarization | 14 |
| 6.1 | Extending summarization | 14 |
| 6.2 | Summarization versus saturation | 15 |
| 6.3 | Shortcut results | 17 |
| 6.4 | Relationships between summaries | 19 |
| 7 | Summarization algorithms | 20 |
| 7.1 | Data graph summarization | 20 |
| 7.2 | Typed graph summarization | 24 |
| 8 | Experimental study | 24 |
| 9 | Related Work | 27 |
| 10 | Conclusion | 28 |
| A | Proof of Proposition 1 | 31 |
| B | Proof of Proposition 2 | 31 |
| C | Proof of Proposition 3 | 31 |
| D | Proof of Theorem 1 | 31 |
| E | Proof of Theorem 2 | 32 |
| F | Proof of Lemma 1 | 34 |
| G | Proof of Lemma 2 | 34 |
| H | Proof of Proposition 4 | 35 |
| I | Proof of Theorem 3 | 35 |
| J | Proof of Lemma 3 | 37 |

| | |
|---------------------------------|-----------|
| K Proof of Proposition 5 | 37 |
| L Proof of Theorem 4 | 37 |
| M Proof of Theorem 7 | 37 |
| N Proof of Proposition 6 | 38 |
| O Proof of Proposition 7 | 40 |

1 Introduction

Large, complex data graphs are everywhere, from social networks to scientific data, biology and bibliography databases etc. *Data graph summaries*, small-size graphs extracted from the data, have been extensively investigated, with two main goals: (i) Speed up graph query processing by deriving an *index* from the summary; (ii) Use the summary *instead* of the graph, for instance, to *allow humans to explore the structure of a new graph through a GUI*, or to decide that a query *lacks answers on the graph* (thus, save the fruitless effort of evaluating it and instead return the empty result without consulting the graph). Numerous graph summaries have been proposed in the past, e.g., [15, 11] (see also Section 9). Among these, many are *graph quotients*: based on an *equivalence notion* \equiv among graph nodes, a quotient summary has a node for each equivalence class. Quotient summaries are easy to interpret by human users, and have interesting properties from a computational perspective (see also Section 3).

The starting point of our work is *quotient summarization of RDF graphs*, the W3C standard for representing Semantic Web data. RDF graphs (and also many non-RDF ones) have a high degree of *structural heterogeneity*: graph nodes may be sources and/or targets of very different sets of properties. Further, RDF allows attaching zero, one or more *types* to nodes; this requires applications to exploit type information when available, but also use graphs partially or completely untyped. Finally, RDF graphs may also include an *ontology* describing relationships between the different node types and edges labels (properties) present in the graph. Ontologies capture crucial domain knowledge about an application, and they allow to *interpret and enrich* the graph based on such knowledge. For instance, an ontology stating that any node with an ISBN number is of type *publication*, on a graph stating that p has the ISBN i , allows to *infer* that p is a publication. Such *implicit* information crucially contributes to a graph's *semantics*; for instance, if a query q requires “all publications with an ISBN”, q would be empty without taking inference into account, whereas the correct answer as defined by the SPARQL query language for RDF graphs should include p .

In the sequel of the paper, Section 2 introduces the graphs we consider, then Section 3 recalls quotient summaries.

Next, we present our contributions:

(i) We propose *two new structural node equivalence relations* (applicable to any labeled, directed graph), based on a novel concept of *data property cliques* (Section 4). Different from known node equivalence relations, which are based on *conjunctions* of structural conditions satisfied by equivalent nodes, our equivalences are based on *disjunctive* conditions. While this makes our summaries less suitable to serve as indexes for conjunctive queries, they are *very tolerant of structural heterogeneity*, thus they are: typically very *compact* (orders of magnitude smaller than the graph), and very suitable for *summary-based graph exploration and visualization*. Indeed, their compactness was an argument of choosing them for integration into LODAtlas [32], an interactive portal for navigating and searching through Linked Open Data.

(ii) When graphs have type edges, we identify two methods for extending *any* node equivalence relation and accordingly, extend *any quotient summarization method* to such edges. We derive concrete

| | Any data graphs | RDF data graphs |
|------------------------|-------------------|-------------------------|
| Any \equiv relation | – | Sections 6.1 and 6.2 |
| Our \equiv relations | Section 4 and 7.1 | Sections 5, 6.3 and 7.2 |

Table 1: Outline of the applicability of our contributions.

summaries for such graphs based on our novel equivalence relations (Section 5).

(iii) When graphs have ontologies, we provide the first *in-depth analysis of the intricate interplay between summarization and saturation*, the reasoning mechanism used to complement an RDF graph with the implicit data it contains (Section 6). Specifically: (iii.a) We show that these two operations do not always commute. Next, we identify a *sufficient condition on a node equivalence relation* which enables summarizing *the explicit and implicit* edges of a graph, without materializing the implicit ones, through a method we call *shortcut*. This method may be much faster than the one which infers all implicit edges before summarizing. (iii.b) We formally establish that two of our equivalence relations allow the shortcut while two others do not.

(iv) We provide *efficient (linear-time in the size of the graph) algorithms* for building our summaries, which are also *incremental*: an edge added to (or removed from) the graph can be reflected in *constant time* in the summary (under an assumption frequently verified by real-life graphs) (Section 7).

Finally, (v) we experimentally validate the good properties of our algorithms (Section 8). Table 1 recaps the areas in which our contributions apply; the code is available online at <https://team.inria.fr/cedar/projects/rdfsummary/>. We discuss related work in Section 9 then we conclude.

Proofs of our technical results are available in the Appendix.

2 Data graphs

Our work is targeted to *directed graphs, with labeled nodes and edges*. This includes those described in RDF [37], the W3C standard for representing Web data. However, RDF graphs attach special interpretation to certain kinds of edges: (i) *type* edges may be used to attach type information to a data node; (ii) *ontology* edges may describe application-domain knowledge as relationships that hold between edge labels and/or node types.

Below, we introduce the useful terminology for RDF graphs since they are *the most general class of graphs for which our summarization methods apply*. We also isolate significant subsets of such graphs, which will be handled differently during summarization.

An *RDF graph* is a set of *triples* of the form $s p o$. A triple states that its *subject* s has the *property* p , and the value of that property is the *object* o . We consider only well-formed triples, as per the RDF specification [37], using uniform resource identifiers (URIs), typed or untyped literals (constants) and blank nodes (unknown URIs or literals). Blank nodes are essential features of RDF allowing to support *unknown URI/literal tokens*. These are conceptually similar to the labeled nulls or variables used in incomplete relational databases [1], as shown in [14].

Notations. We use s , p , and o as placeholders for subjects, properties and objects, respectively.

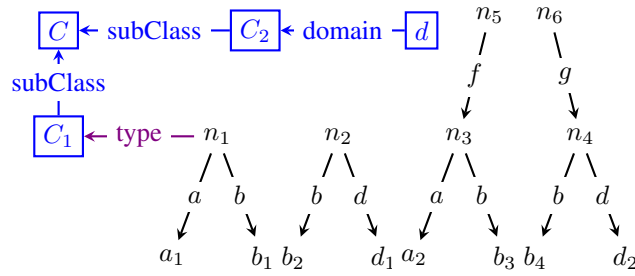
Figure 1 (top) shows how to use triples to describe resources, that is, to express class (unary relation) and property (binary relation) assertions. The RDF standard [37] has a set of *built-in classes and properties*, as part of the `rdf:` and `rdfs:` pre-defined namespaces. We use these namespaces exactly for these classes and properties, e.g., `rdf:type` specifies the class(es) to which a resource belongs.

As our running example, Figure 2 shows a sample RDF graph. Black and violet edges encode data and type respectively, e.g., node n_1 has property a whose object (value) is a_1 , n_1 is of class (or has type)

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

| Constraint | Triple | OWA interpretation |
|---------------|----------------------------|--------------------------------------|
| Subclass | $s \text{ subclass } o$ | $s \subseteq o$ |
| Subproperty | $s \text{ subproperty } o$ | $s \subseteq o$ |
| Domain typing | $p \text{ domain } o$ | $\Pi_{\text{domain}}(s) \subseteq o$ |
| Range typing | $p \text{ range } o$ | $\Pi_{\text{range}}(s) \subseteq o$ |

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

Figure 2: Sample RDF graph $G = \langle D_G, S_G, T_G \rangle$: D_G edges are shown in black, S_G edges in blue, T_G edges in violet.

C_1 , etc.

RDF Schema (RDFS). RDFS allows enhancing the descriptions in RDF graphs by declaring *ontological constraints* between the classes and the properties they use. Figure 1 (bottom) shows the four kinds of RDFS constraints, and how to express them through triples. Here, “domain” denotes the first, and “range” the second attribute of every property. In Figure 2, the blue edges connecting boxed nodes are RDFS constraints: they state that C_1 and C_2 are subclasses of C , and that the domain of d is C_2 .

RDFS constraints are interpreted under the *open-world assumption (OWA)* [1], i.e., as deductive constraints. *RDF entailment* is the mechanism through which, based on a set of explicit triples and some *entailment rules*, implicit RDF triples are derived. For instance, in Figure 2, node n_1 is of type C_1 , which is a subclass of C . Through RDF entailment with the subclass constraint in Figure 1, we obtain the *implicit (entailed) triple* $n_1 \text{ type } C$ stating that n_1 is of class C . Similarly, n_2 and n_4 have property d whose domain is C_2 (thanks to the domain constraint in Figure 1), thus n_2 and n_4 are also of class C_2 . Further, because C_2 is a subclass of C , they are also of class C .

In general, 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 entailment rules that leads from G to $s \text{ p } o$ (where at each step, triples previously entailed are also taken into account).

RDF graph saturation. Given a set of entailment rules, the *saturation* (a.k.a. closure) of an RDF graph G , is defined as the fixpoint obtained by repeatedly adding to G the triples derived using the entailment rules; we denote it G^∞ . For the four constraints shown in Figure 1, which we consider throughout this work, the saturation of G , is finite, 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 entails (leads to, has as logical consequence) the triple $s \text{ p } o$ if and only if $s \text{ p } o \in G^\infty$. It is important to note that the semantics of an RDF graph is its saturation [37].

In particular, when querying an RDF graph, the *answer* to the query should be computed both from its explicit and its implicit triples.

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

1. Triple-based representation of an RDF graph. We see G as a *union of three edge-disjoint subgraphs* $G = \langle D_G, S_G, T_G \rangle$, where: (i) S_G , the **schema** component, is the set of all G triples whose properties are *subclass*, *subproperty*, *domain* or *range*; we depict such triples with **blue** edges; (ii) T_G , the **type** component, is the set of *type* triples from G ; we show them in **violet**; (iii) D_G , the **data** component, holds all the remaining triples of G ; we display them in black. Note that each union of the D_G , S_G , and T_G components 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 .

2. The graph-based representation of an RDF graph. As per the RDF specification [37], the *set of nodes* of an RDF graph is the set of subjects and objects of triples in the graph, while its *edges* correspond to its triples. We define three categories of RDF graph nodes: (i) a *class node* is any node whose URI appears as subject or object of a *subclass* triple, or object of a *domain* or *range* or *type* triple; we show them in blue boxes; (ii) a *property node* is any node whose URI appears as subject or object of a *subproperty* triple, or subject of a *domain* or *range* triple, or property of a triple¹; in Figure 2, d is a property node. We also show them in blue boxes; (iii) a *data node* as any node that is neither a class nor a property node. We show them in black. Note that the sets of class nodes and of property nodes may intersect (indeed, nothing in the RDF specification forbids it). However, data nodes are disjoint from both class and property nodes.

We will rely on the graph, respectively, the triple-based representation when each is most natural for the presentation; accordingly, we may use *triple* or *edge* interchangeably to denote a graph edge.

3 Summarization framework

We recall the classical notion of *graph quotient*, on which many graph summaries, including ours, are based. In particular, we recall quotients based on *bisimilarity*, and show that their very nature makes them ill-suited to summarize heterogeneous graphs.

Graph quotients. Given a data graph G and an equivalence relation² \equiv over the node of G , the quotient of G by \equiv , denoted $G_{/\equiv}$, is the graph having (i) a node for each equivalence class of \equiv (thus, for each set of equivalent G nodes); and (ii) for each edge $n_1 \xrightarrow{a} n_2$ in G , an edge $m_1 \xrightarrow{a} m_2$, where m_1, m_2 are the quotient nodes corresponding to the equivalence classes of n_1, n_2 respectively.

Many known graph summaries, e.g., [29, 20, 9, 35, 24, 11, 13] are quotient-based; they differ in their equivalence relations \equiv . Quotient summaries have several desirable properties:

Size guarantees By definition, $G_{/\equiv}$ is guaranteed to have at most as many nodes and edges as G . Some non-quotient summaries, e.g., Dataguides [15], cannot guarantee this.

Property completeness denotes the fact that every property (edge label) from G is present on summary edges. This is helpful to users approaching a graph dataset for the first time³; it is for this scenario precisely that our summaries are used in LODAtlas [32].

¹A property node must be a *node*, i.e., merely appearing in a property position does not make an URI a property node; for this, the URI needs to appear *as a subject or object* in a triple of the graph.

²An equivalence relation \equiv is a binary relation that is reflexive, i.e., $x \equiv x$, symmetric, i.e., $x \equiv y \Rightarrow y \equiv x$, and transitive, i.e., $x \equiv y$ and $y \equiv z$ implies $x \equiv z$ for any x, y, z .

³Most complete summarization methods, including ours, can be adapted to reflect e.g., only properties above a certain frequency threshold in the graph etc. We will not consider this further.

Structural representativeness is the following property: for any query q that has answers on G , its *structure-only* version q' , which copies all the graph patterns of q but erases its possible selections on nodes as well as counting predicates, is guaranteed to have answers on $G_{/\equiv}$.

For instance, if q_1 is “find all nodes that are target of an f edge and source of a b and a d edge” on the graph in Figure 2, then q'_1 is the same as q_1 . If the query q_2 is “find all nodes whose labels contain “Alice”, having exactly one (not more) incoming f edge and exactly one outgoing b edge”, the query q'_2 asks for “all nodes having incoming f and outgoing b edges”. Thanks to representativeness, quotient summaries can be used to prune empty-answer queries: if $q'(G_{/\equiv})$ has no answers, then q has no answers on G . Since the summary is often much smaller than G , pruning is very fast and saves useless query evaluation effort on G .

To enjoy the above advantages, in this work, *a summary of G is its quotient through some equivalence relation*.

Two extreme quotient summaries help to see the trade-offs in this context. First, let \top denote the equivalence relation for which all nodes are equivalent: then, $G_{/\top}$ has a single node with a loop edge to itself for each distinct property in G . This summary collapses (and loses) most of the graph structure. Now, let \perp denote the equivalence relation for which each node is only equivalent to itself. Then, $G_{/\perp}$ is isomorphic to G for any graph G ; it preserves all the structure but achieves no summarization.

Bisimulation-based summaries. Many known structural quotient summaries, e.g., [29, 9, 21, 12] are based on bisimilarity [19]. Two nodes n_1, n_2 are *forward and/or backward bisimilar* (denoted \equiv_{fw} , \equiv_{bw} and \equiv_{fb}) iff for every G edge $n_1 \xrightarrow{a} m_1$, G also comprises an edge $n_2 \xrightarrow{a} m_2$, such that m_1 and m_2 are also forward and/or backward bisimilar, respectively. The \equiv_{fw} and \equiv_{bw} relations only take into account the paths outgoing from (resp. only the paths incoming to) the nodes. The **symmetry** of $G_{/\equiv_{fb}}$ is an advantage, as it makes it more resistant to minor modeling differences in the data. For instance, let t' be the triple a hasAuthored p and t'' be p hasAuthor a , which essentially represent the same information. Triple t' would impact the nodes to which a is \equiv_{fw} , while it would not impact the nodes to which a is \equiv_{bw} ; symmetrically, t'' would impact the \equiv_{bw} class of p but not its \equiv_{fw} class. In contrast, \equiv_{fb} reflects this information whether it is modeled in one direction or in the other.

We denote the bisimulation based summaries $G_{/\equiv_{fw}}$ (forward), $G_{/\equiv_{bw}}$ (backward) and $G_{/\equiv_{fb}}$ (forward and backward), respectively. They tend to be **large** because *bisimilarity is rare in heterogeneous data graphs*. For instance, each node of the graph in Figure 2 is only \equiv_{fb} to itself, thus \equiv_{fb} is useless for summarizing it; our experiments in Section 8 confirm this on many graphs. To mediate this problem, **k -bisimilarity** has been introduced [22], whereas nodes are k -forward (and/or backward) bisimilar iff their adjacent paths of length at most k are identical; the smaller k is, the more permissive the equivalence relation.

One drawback of k -bisimilarity is that it requires users to guess the k value leading to the best compromise between compactness (which favors low k ; note that $k = 0$ leads exactly to $G_{/\top}$) and structural information in the summary (high k). Further, *even 1-bisimilarity is hard to achieve* in heterogeneous graphs. For instance, Figure 3 shows the 1fb summary of the sample graph in Figure 2⁴. Nodes in the bottom row of G , which only have incoming a , b , respectively d edges, and have no outgoing edges, are summarized together. However, none of n_1, n_2, n_3 and n_4 are equivalent, because of the presence/absence of a and d edges, and of their possible incoming edges.

Any structural equivalence relation cuts a trade-off between compactness and structure preservation. Below, we introduce two relations leading to compact summaries that *in addition* cope well with the data heterogeneity frequently encountered in RDF (and other) graphs.

Other interesting properties of such relations are: the complexity of building the corresponding summary, and of updating it to reflect graph updates. Further, the presence of type and schema (ontology) triples has not been formally studied in quotient graph summarization.

⁴In Figure 3, type and schema triples are summarized according to the method we propose below. However, the treatment of these triples is orthogonal to the interesting aspects of this example.

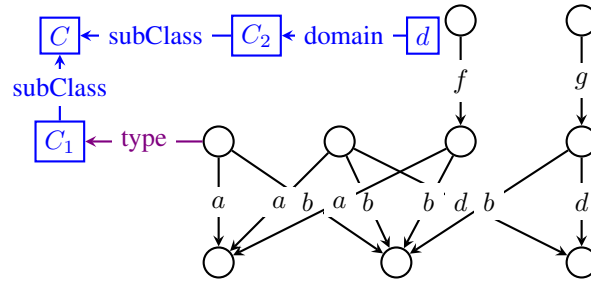


Figure 3: 1fb summary of the RDF graph in Figure 2.

4 Data graph summarization

We first consider graphs made of *data triples* only, thus of the form $G = \langle D_G, \emptyset, \emptyset \rangle$. We define the novel notion of *property cliques* in Section 4.1; building on them, we devise new graph node equivalence relations and corresponding graph summaries in Section 4.2. Summarization will be generalized to handle also *type triples* in Section 5, and *type and schema triples* in Section 6.

4.1 Data property cliques

We are interested in defining equivalence relations between two data nodes, that can cope with the heterogeneity present in many real-life data graphs. For example, up to 14 data properties (such as title, author, year, but also note, etc.) are used to describe conference papers in a graph version of the DBLP bibliographic database. Each paper has a certain subset of these 14 properties, and has some of them, e.g., authors, with multiple values; we counted more than 130 such distinct property subsets in a small (8MB) fragment of DBLP. To avoid the “noise” introduced by such structural heterogeneity, we need node equivalence relations that look *beyond* it, and consider that all the nodes corresponding to conference publications are equivalent.

To do that, we first focus on the way data properties (edge labels) are organized in the graph. The simplest relation that may exist between two properties is *co-occurrence*, when a node is the source (or target) of two edges carrying the two labels. However, in heterogeneous RDF graphs such as DBLP, two properties, say author and title, may co-occur on a node n , while another node n' has title, year, and howpublished: we may consider *all* these properties (author, title, year and howpublished) related, as they (directly or *transitively*) co-occur on some nodes. Formally:

Definition 1. (PROPERTY RELATIONS AND CLIQUES) *Let p_1, p_2 be two data properties in G :*

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

In the graph in Figure 2, properties a and b are source-related due to n_1 (condition 1. in the definition). Similarly, b and d are source-related due to n_2 ; consequently, a and d are source-related (condition 2.).

| n | n_1 | n_2 | n_3 | n_4 | n_5 | n_6 |
|---------|---------------|---------------|---------------|---------------|-------------|-------------|
| $SC(n)$ | $\{a, b, d\}$ | $\{a, b, d\}$ | $\{a, b, d\}$ | $\{a, b, d\}$ | $\{f\}$ | $\{g\}$ |
| $TC(n)$ | \emptyset | \emptyset | $\{f\}$ | $\{g\}$ | \emptyset | \emptyset |

| n | a_1 | a_2 | b_1 | b_2 | b_3 | d_1 | d_2 |
|---------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| $SC(n)$ | \emptyset | \emptyset | \emptyset | \emptyset | \emptyset | \emptyset | \emptyset |
| $TC(n)$ | $\{a\}$ | $\{a\}$ | $\{b\}$ | $\{b\}$ | $\{b\}$ | $\{d\}$ | $\{d\}$ |

Table 2: Source and target cliques of G nodes (Figure 2).

Thus, a source clique of this graph is $SC_1 = \{a, b, d\}$. Table 2 shows the target and source cliques of all data nodes from Figure 2.

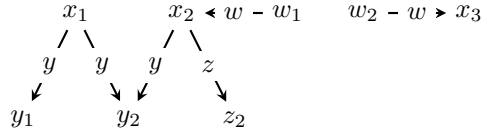
It is easy to see that the set of non-empty source (or target) property cliques is a *partition over the data properties of G*. Further, if a node $n \in G$ is source of some data properties, they are all in the same source clique; similarly, all the properties of which n is a target are in the same target clique. This allows us to refer to *the source (or target) clique of n* , denoted $SC(n)$ and $TC(n)$.

4.2 Strong and weak node equivalences

Building on property cliques, we define two main *node equivalence relations* among the data nodes of a graph G:

Definition 2. (STRONG EQUIVALENCE) *Two data nodes of G are strongly equivalent, denoted $n_1 \equiv_S n_2$, iff they have the same source and target cliques.*

Strongly equivalent nodes have the same structure of *incoming and outgoing* edges. In Figure 2, nodes n_1, n_2 are *strongly* equivalent to each other, and so are e.g., a_1, a_2, b_1, b_2 and b_3 etc.

Figure 4: Sample weakly equivalent nodes: x_1, x_2, x_3 .

A second, weaker notion of node equivalence could request only that equivalent nodes share the same *incoming or outgoing* structure, i.e., they share the same source clique or the same target clique. Figure 4 illustrates this. Nodes x_1, x_2 have the same source clique because they both have outgoing y edges. Further, x_2 and x_3 have the same target clique because both have incoming w edges. Since equivalence must be transitive, it follows that x_1 and x_3 must also be considered weakly equivalent, since they “follow the same pattern” of having at least one incoming w edge, or at least one outgoing y or z edge, and no other kinds of edges. Formally:

Definition 3. (WEAK EQUIVALENCE) *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.*

It is easy to see that \equiv_W and \equiv_S are equivalence relations and that strong equivalence implies weak equivalence.

In Figure 2, n_1, \dots, n_4 are *weakly* equivalent to each other due to their common source clique SC_1 ; a_1, a_2 are weakly equivalent due to their common target clique etc.

4.3 Weak and strong summarization

Notation: representation function. We say the summary node of G/\equiv corresponding to the equivalence class of a G node n represents n , and denote it $f_{\equiv}(n)$ or simply $f(n)$ when this does not cause confusion. We call f_{\equiv} the *representation function* of the equivalence relation \equiv over G .

Weak summarization. The first summary we define is based on weak equivalence:

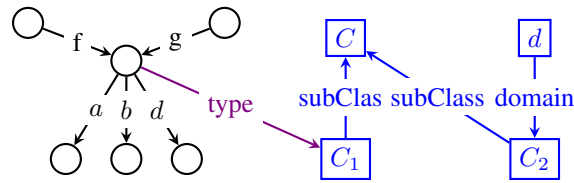


Figure 5: Weak summary of the RDF graph in Figure 2.

Definition 4. (WEAK SUMMARY) *The weak summary of a data graph G , denoted G/w , is its quotient graph w.r.t. the weak equivalence relation \equiv_w .*

The weak summary of the graph in Figure 2 is depicted by the black nodes and edges in Figure 5⁵; summary nodes are shown as unlabeled circles, to denote that they are anonymous (new) nodes, each of which represents one or more G nodes. The central one represents n_1, n_2, n_3 and n_4 . Its outgoing edges go towards nodes representing, respectively, a_1 and a_2 ; b_1, b_2 and b_3 ; finally, d_1 and d_2 . Its incoming edges come from the representative of n_5 (which was a source of an f edge in G), respectively from the representative of n_6 (source of g).

The weak summary has the following important property:

Proposition 1. (UNIQUE DATA PROPERTIES) *Each G data property appears exactly once in G/w .*

Importantly, the above Proposition 1 warrants that $|G/w|$, the number of edges in G/w , is exactly the number of distinct data properties in G . This observation is used in our weak summarization algorithms (Section 7). By definition of a quotient summary (Section 3), *this is the smallest number of edges a summary may have* (since it has at least one edge per each distinct property in G). Thus, G/w is a minimal-size quotient summary (like G/\top from Section 3, but much more informative than it). As our experiments show, $|G/w|$ is typically 3 to 6 orders of magnitude smaller than $|G|$.

Strong summarization. Next, we introduce:

Definition 5. (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 .*

The strong summary of the graph of Figure 2 is shown by the black edges and nodes in Figure 6. Similarly to the weak summary (Figure 5), the strong one features a single node source of a , respectively, b, d, f and g edges. However, differently from G/w , the strong summary splits the data nodes whose source clique is $\{a, b, d\}$ in *three equivalence classes*: n_1 and n_2 have the empty target clique, while that of n_3 is $\{f\}$ and that of n_4 is $\{g\}$. Thus, two data nodes represented by the same strong summary node have similar structure both in their inputs and outputs; in contrast, a weak summary (recall Figure 4)

⁵The violet and blue edges serve our discussion later on.

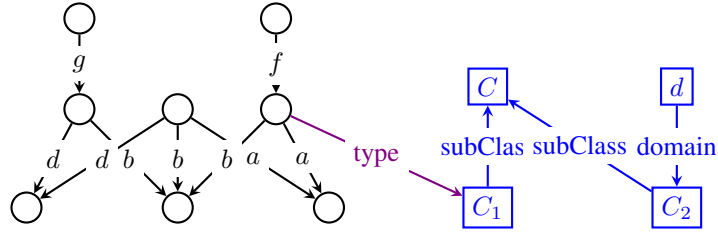


Figure 6: Strong summary of the RDF graph in Figure 2.

represents together nodes having similar structure in their inputs *or* outputs, *or* which are both equivalent to another common node. As we can see, *strong summarization leads to finer-granularity* summaries. An effect of this finer granularity is that in $G_{/s}$, several edges may have the same label, e.g., there are three edges labeled b in Figure 5 (whereas for $G_{/w}$, as stated in Proposition 1, this is not possible). Our experiments (Section 8) show that while $G_{/s}$ is often somehow larger than $G_{/w}$, it still remains many orders of magnitude smaller than the original graph.

By definition of \equiv_s , equivalent nodes have the same source clique and the same target clique. This leads to:

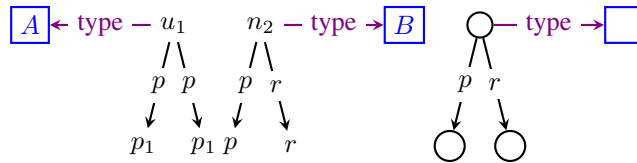
Proposition 2. (STRONG SUMMARY NODES AND G CLIQUES) $G_{/s}$ has exactly one node for each source clique and target clique of a same node $n \in \mathcal{D}_G$.

Proposition 2 is exploited by the implementations of our strong summarization algorithms (Section 7).

5 Typed data graph summarization

We generalize our approach to summarize graphs with data and type triples, thus of the form $G = \langle \mathcal{D}_G, \mathcal{T}_G, \emptyset \rangle$.

Starting from an equivalence relation \equiv defined over *data* nodes, in order to summarize $\mathcal{D}_G \cup \mathcal{T}_G$, two questions must be answered: (i) how should \equiv be extended on class nodes (such as C_1 in Figure 2)? and (ii) how should the type edge(s) of a node n be taken into account when determining to whom n is equivalent? Below, we answer these questions for *any* equivalence relation \equiv , then instantiate our answer to the weak and strong relations we defined.



To study the first question, consider the sample graph above, and a possible summary of this graph at its right. Assume that the types A and B are considered equivalent. Quotient summarization represents them both by the summary node at the top right, which (like all summary nodes) is a “new” node, i.e., it is neither A nor B . Observe that this summary *compromises representativeness* for queries over both the data and the type triples: for instance, the query asking for “nodes of type A having property r ” is empty on the summary (as type A has been conflated with type B) while it is non empty on the graph.

To avoid this, we argue that when moving from data to typed data graphs, any equivalence relation \equiv between data nodes should be extended to class nodes as follows: 1. *any class node is only equivalent*

to itself and 2. any class node is only represented by itself, hence a graph has the same class nodes as its summary.

We now consider the second question: how should \equiv be extended to exploit not only the data but also the type triples? Note that two nodes may have similar incoming/outgoing data edges but different type edges, or vice-versa, the same types but very different data edges. We introduce two main alternatives below, then decline them for weak and strong summarization.

5.1 Data-then-type summarization

This approach consists of using an equivalence \equiv defined based on data properties in order to determine which data nodes are equivalent and thus to build summary nodes, and data edges between them. Afterward, for each triple n type C in G , add to G_{\equiv} a triple $f_{\equiv}(n)$ type C , where we recall that $f_{\equiv}(n)$ is the representative of n in G_{\equiv} . This approach is interesting, e.g., when only some of the nodes have types (often the case in RDF graphs). In such cases, it makes sense to first group nodes according to their data edges, while still preserving the (partial) type information they have. We extend the W , respectively S summaries to type triples, by stating that they (i) represent each class node by itself; and (ii) follow a data-then-type approach, as described above.

In Figure 5, the black and violet edges (including the C_1 node) depict the weak summary of the black and violet graph triples Figure 2. The type edge reads as: *at least one* of the nodes represented by its source, was declared of type C_1 in the input graph. Similarly, the black and violet edges in Figure 6 show the strong summary of the same subset of our sample graph.

To recap, in data-then-type summarization using \equiv , two data nodes are represented together iff they are \equiv based on their incoming and outgoing data edges, while a class node is only equivalent to itself (and always represented by itself).

One more point needs to be settled. Some T_G nodes may have types, but no incoming or outgoing data properties. Strong summarization represents all such nodes together, based on their (\emptyset, \emptyset) pair of source and target cliques. For completeness, we extend weak summaries to also represent such nodes together, by a single special node denoted N_{\emptyset} .

5.2 Type-then-data summarization

This approach takes the opposite view that node types are more important when deciding whether nodes are equivalent. Observe that our framework (just like RDF) does not prevent a node from having *several* types. At the same time, representing a node by *each* of its types separately would violate the quotient summarization framework, because a quotient, by definition, represents each node exactly once. Thus, in type-then-data summarization, we extend a given equivalence relation \equiv (based on data properties alone) as follows.

Definition 6. (TYPED EQUIVALENCE) Typed equivalence, denoted \equiv_T , is an equivalence relation over $D_G \cup T_G$ 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; any class node is only equivalent to itself.

Intuitively, typed equivalence performs a first-cut data node classification, according to their sets of types. In particular, all untyped nodes are equivalent to themselves. This enables the definition of type-then-data summaries as *double quotients*: first, quotient G by \equiv_T ; then, quotient the resulting graph by some data node equivalence *only on untyped nodes* (each left alone in an equivalence class of \equiv_T), to group them according to their data edges.

Applied to weak summarization, this approach leads to:

Definition 7. (TYPED WEAK SUMMARY) Let \equiv_{UW} (untyped weak equivalence) be an equivalence relation that holds between two data 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 a graph G is the untyped-weak summary of the type-based summary of G , namely $(G_{/T})_{/UW}$.

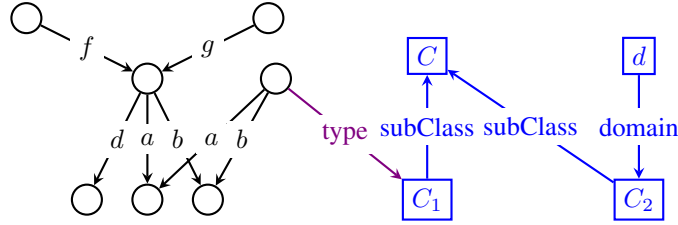


Figure 7: Typed weak summary of the graph in Figure 2.

In Figure 7, the black and violet edges depict the typed weak summary of the data and type edges of the sample graph in Figure 2. Unlike $G_{/W}$ (Figure 5), $G_{/TW}$ represents the node of type C_1 separately from the untyped ones having similar properties. This reflects the primordial role of types in type-then-data summarization.

In a similar manner, we define:

Definition 8. (TYPED STRONG SUMMARY) Let \equiv_{US} (untyped strong equivalence) be an equivalence relation that holds between two data 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: $(G_{/T})_{/US}$.

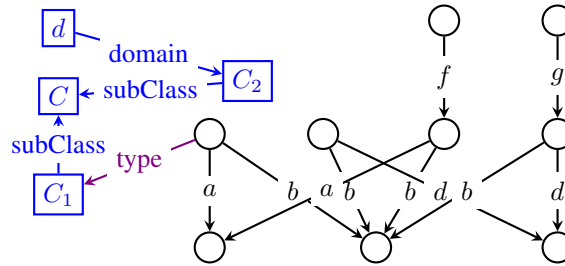


Figure 8: Typed strong summary of the graph in Figure 2.

In Figure 8, the black and violet edges depict the typed strong summary of the data and type edges of the sample graph in Figure 2. Unlike $G_{/S}$ (Figure 6), $G_{/TS}$ represents the node n_1 of type C_1 separately from n_2 , which has no types.

6 RDF graph summarization

We consider now the summarization of general RDF graphs which may also have schema triples, i.e., of the form $G = \langle D_G, T_G, S_G \rangle$.

6.1 Extending summarization

First, how to extend an equivalence relation \equiv defined on data nodes (and extended as we discussed in Section 5 to class nodes, each of which is only equivalent to itself), to also cover *property nodes*, such as the boxed d node in Figure 2? Such property nodes provide important schema (ontology) information,

which describes how the properties and types present in the data relate to each other, and leads to implicit triples (recall Section 2). For the summary to preserve the semantics of the input graph, by an argument similar to the one at the beginning of Section 5, we impose that \equiv be extended so that *any property node should (also) be equivalent only to itself*, and propose that in any summary, *any property node should be represented by itself*. As a consequence, since any class or schema node is equivalent only to itself and represented only by itself:

Proposition 3. (SCHEMA PRESERVATION THROUGH SUMMARIZATION) *For any equivalent relation \equiv defined on data nodes and extended as specified above to class and property nodes, and any graph $G = \langle D_G, T_G, S_G \rangle$, it holds that: $S_{G_{\equiv}} = S_G$, that is: the summary of G through \equiv has exactly the schema triples of G .*

This decision allows us to simply copy schema triples from the input graph to each of its summaries. Figures 5, 6, 7 and 8, considered in their entirety, show respectively full $G_{/W}$, $G_{/S}$, $G_{/TW}$ and $G_{/TS}$ summaries of the sample RDF graph in Figure 2.

6.2 Summarization versus saturation

As we explained in Section 2, the semantics of a graph G includes its explicit triples, but also its *implicit* triples which are not in G , but hold in G^∞ due to ontological triples (such as the triples n_2 type C_2 and n_2 type C discussed in Section 2).

A first interesting question, then, is: how does saturation impact the summary of a graph? As Figure 1 shows, saturation adds data and type triples. Other entailment rules (see [14]) also generate schema triples, e.g., if C' is a subclass of C'' and C'' is a subclass of C''' , then C' is also a subclass of C''' etc. Due to these extra edges, in general, $(G^\infty)_{/\equiv}$ and $G_{/\equiv}$ are different. On one hand, their nodes may be different, but this is not the most interesting aspect, as summary nodes are just “representatives”, i.e., they are labeled in a somehow arbitrary fashion. On the other hand (and this is much more meaningful), their graph structure may be different, as it results from graphs with different edges. To separate the mere difference of node IDs from the meaningful difference of graph structure, we define:

Definition 9. (STRONG ISOMORPHISM \simeq) *A strong isomorphism between two RDF graphs G_1, G_2 , noted $G_1 \simeq G_2$, is a graph isomorphism which is the identity for the class and property nodes.*

Intuitively, strongly isomorphic graphs (in particular, summaries) represent exactly the same information, while the identifiers of their non-class, non-property nodes (shown as unlabeled circles in our examples) may differ.

Next, one could wonder whether saturation *commutes with* summarization, that is, does $(G^\infty)_{/\equiv} \simeq (G_{/\equiv})^\infty$ hold? If this was the case, it would lead to a likely more efficient method for computing the summary of G 's full semantics, that is $(G^\infty)_{/\equiv}$, without saturating G (thus, without materializing all its implicit triples); instead, we would summarize G and then saturate the resulting (usually much smaller) graph. Unfortunately, Figure 9 shows that this is not always the case. For a given graph G , it traces (top row) its weak summary $G_{/W}$ and its saturation $(G_{/W})^\infty$, whereas the bottom row shows G^∞ and its summary $(G^\infty)_{/W}$. Here, saturation leads to b edges outgoing both r_1 and r_2 which makes them equivalent. In contrast, summarization *before* saturation represents them separately; saturating the summary cannot revert this decision, to unify them as in $(G^\infty)_{/W}$ (recall from Section 2 that saturation can only add edges between G nodes).

While $(G^\infty)_{/\equiv}$ and $(G_{/\equiv})^\infty$ are *not* strongly isomorphic in general, we establish that they *always* relate as follows:

Theorem 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. Then f defines a homomorphism from G^∞ to $(G_{/\equiv})^\infty$.*

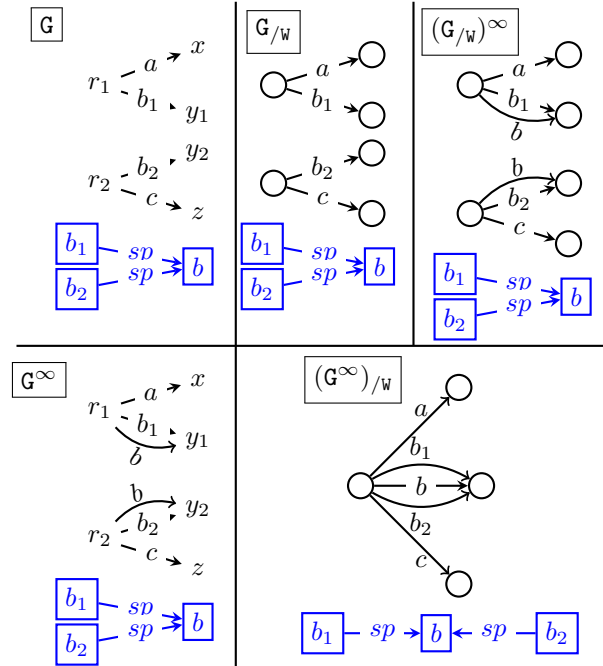


Figure 9: Saturation and summarization on a sample graph.

Since we established that $(G/\equiv)^\infty$ is homomorphic to G^∞ , would *their* summaries be the same? In other words: are $(G^\infty)_{/\equiv}$ and $((G/\equiv)^\infty)_{/\equiv}$ strongly isomorphic?

It turns out that this may hold or not depending on the RDF equivalence relation under consideration. When it holds, we call *shortcut* the following three-step transformation aiming at obtaining $(G^\infty)_{/\equiv}$: first summarize G ; then saturate its summary; finally, summarize it again in order to build $((G/\equiv)^\infty)_{/\equiv}$:

Definition 10. (SHORTCUT) *We say the shortcut holds for a given RDF node equivalence relation \equiv iff for any RDF graph G , $(G^\infty)_{/\equiv}$ and $((G/\equiv)^\infty)_{/\equiv}$ are strongly isomorphic.*

Next, we establish one of our main contributions: a *sufficient condition* under which for any quotient summary based on an equivalence relation \equiv as discussed above (where class and property nodes are preserved by summarization), the shortcut holds. In particular, as we will demonstrate (Section 8), the existence of the shortcut can lead to computing $(G^\infty)_{/\equiv}$ *substantially faster*.

Theorem 2 (Sufficient shortcut condition). *Let $G_{/\equiv}$ be a summary of G through \equiv and $f_{/\equiv}$ the corresponding representation function from G nodes to $G_{/\equiv}$ nodes (see Figure 10).*

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

Figure 10 depicts the relationships between an RDF graph G , its saturation G^∞ and summarization $(G^\infty)_{/\equiv}$ thereof, and the RDF graphs that appear at each step of the shortcut computation. The **intuition for the sufficient condition** is the following. On any path in Figure 10, saturation adds edges to its input graph, while summarization “fuses” nodes into common representatives. On the regular path from G to $(G^\infty)_{/\equiv}$, edges are added in the first step, and nodes are fused in the second. On the shortcut (green) path, edges are added in the second step, while nodes are fused in the first and third steps. The two paths starting from G can reach \simeq results only if G nodes fused on the shortcut path, are also fused (when summarizing

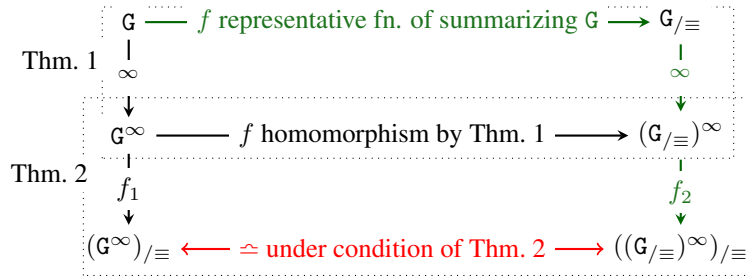


Figure 10: Illustration for Theorem 1 and Theorem 2.

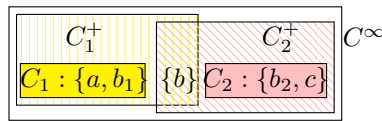


Figure 11: Two source cliques from the graph in Figure 9, their saturations, and their enclosing clique in G^∞ .

G^∞) on the standard path. In particular, *the first summarization along the shortcut path should not make wrong node fusions*, that is, fusions not made when considering the full G^∞ : such a “hasty” fusion *can never be corrected later on along the shortcut path*, as neither summarization nor saturation split nodes. Thus, an erroneous fusion made in the first summarization step irreversibly prevents the end of the shortcut path from being \simeq to $(G/\equiv)^\infty$. In this case, we remark that $((G/\equiv)^\infty)/\equiv$ is just homomorphic to G^∞ (through f then f_2 in Figure 10), hence a compression of it.

6.3 Shortcut results

Next, we investigate for which \equiv does the shortcut hold.

The impact of saturation on property cliques is at the core of our findings. In G^∞ , every G node has all the data properties it had in G , therefore two data properties belonging to a G clique are also in the same clique of G^∞ . Further, if the schema of G comprises *subproperty* constraints, a node may have in G^∞ a data property that it did not have in G . In turn, this leads to G^∞ cliques which “subsume and fuse” several G cliques into one.

For a given clique C of G , we call *saturated clique* and denote C^+ the set 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 nodes in G^∞ , and thus in general C^+ is *not* a clique of G^∞ .

The following Lemma characterizes the relationships between a clique C of G , its saturated version C^+ , and the cliques of G^∞ :

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

1. *There exists exactly one source (resp. target) clique C^∞ of G^∞ such that $C \subseteq C^\infty$.*
2. *If $C_1^+ \cap C_2^+ \neq \emptyset$, then all the properties in C_1 and C_2 are in the same G^∞ clique C^∞ .*
3. *Any non-empty source (or target) clique C^∞ is a union of the form $C_1^+ \cup \dots \cup C_k^+$ for some $k \geq 1$, where each C_i is a non-empty source (resp. target) clique of G , and 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 \mathbb{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^∞ of \mathbb{G}^∞ if and only if there exist k non-empty source (resp. target) cliques of \mathbb{G} , $k \geq 0$, denoted D_1, \dots, D_k such that:

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

Figure 11 illustrates the lemma for two source cliques of the graph in Figure 9. The source clique $C^\infty = \{a, b_1, b, b_2, c\}$ of the \mathbb{G}^∞ (also shown in Figure 9) encloses the cliques C_1, C_2 of \mathbb{G} (Lemma item 1). C_2^+ intersects C_1^+ , thus they are all in the same clique C^∞ (item 2), which is the union of C_1^+ and C_2^+ (item 3). A chain of intersecting cliques connects C_1^+ to C_2^+ . Item 4 follows quite directly from item 3; in Figure 11, it is illustrated e.g., for properties a and c by taking $D_1 = C_1, D_2 = C_2$.

Based on Theorem 2, we show:

Theorem 3 (W shortcut). *The shortcut holds for \equiv_W .*

For instance, on the graph in Figure 9, it is easy to check that applying summarization on $(\mathbb{G}/_W)^\infty$ (as prescribed by the shortcut) leads exactly to a graph strongly isomorphic to $(\mathbb{G}^\infty)/_W$.

Showing Theorem 3 is rather involved. We do it in several steps. First, based on Lemma 1, we show:

Lemma 2 (Property relatedness in W summaries). *Data properties are target-related (resp. source-related) in $(\mathbb{G}/_W)^\infty$ iff they are target-related (resp. source-related) in \mathbb{G}^∞ .*

Based on the above and on Theorem 1, we establish the next result from which Theorem 3 directly follows:

Proposition 4. (SAME CLIQUES-W) \mathbb{G}^∞ and $(\mathbb{G}/_W)^\infty$ have identical source clique sets, and identical target clique sets. Further, a node $n \in \mathbb{G}^\infty$ has exactly the same source and target clique as $f_W(n)$ in $(\mathbb{G}/_W)^\infty$.

Theorem 4 (S shortcut). *The shortcut holds for \equiv_S .*

We prove this based on counterparts of statements established for \mathbb{G}_W . First we show:

Lemma 3 (Property relatedness in S summaries). *Data properties are target-related (resp. source-related) in $(\mathbb{G}/_S)^\infty$ iff they are target-related (resp. source-related) in \mathbb{G}^∞ .*

Then, from Theorem 1 and the above Lemma, we obtain the next proposition from which Theorem 4 directly follows:

Proposition 5. (SAME CLIQUES-S) \mathbb{G}^∞ and $(\mathbb{G}/_S)^\infty$ have identical source clique sets, and identical target clique sets. Further, a node $n \in (\mathbb{G}/_S)^\infty$ has exactly the same source and target clique as $f_S(n)$ in $(\mathbb{G}/_S)^\infty$.

Finally, we have:

Theorem 5 (No shortcut for \equiv_{TW}). *The shortcut does not hold for \equiv_{TW} .*

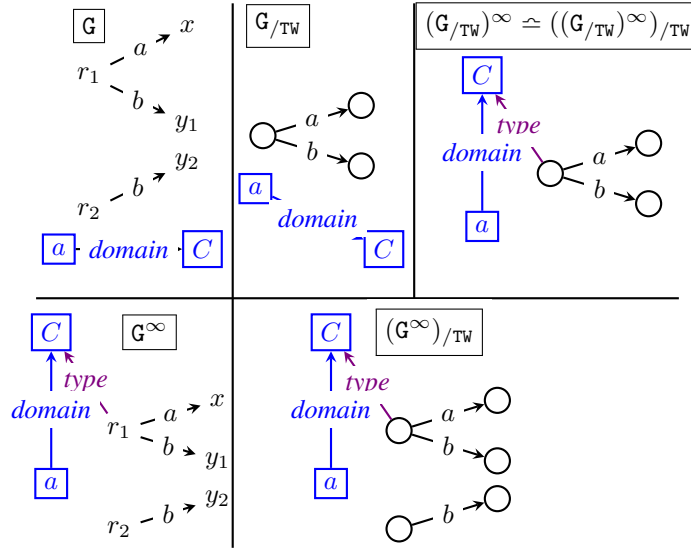


Figure 12: Shortcut counter-example.

We prove this by exhibiting in Figure 12 a counter-example. In G and $G_{/TW}$, all data nodes are untyped; only after saturation a node gains the type C . Thus, in $G_{/TW}$, one (untyped) node represents all data property subjects; this is exactly a “hasty fusion” as discussed below Theorem 2. In $(G_{/TW})^\infty$, this node gains a type, and in $((G_{/TW})^\infty)_{/TW}$, it is represented by a single node. In contrast, in G^∞ , r_1 is typed and r_2 isn't, leading to two distinct nodes in $(G^\infty)_{/TW}$. This is not isomorphic with $(G_{/TW})^\infty$ which, in this example, is strongly isomorphic to $((G_{/TW})^\infty)_{/TW}$. Thus, the shortcut does not hold for $\equiv_{/TW}$ does not admit a shortcut.

Theorem 6 (No shortcut for \equiv_{TS}). *The shortcut does not hold for \equiv_{TS} .*

The graph in Figure 12 is also a shortcut counter-example for TS. More generally, let \equiv_X be an arbitrary RDF node equivalence, and \equiv_{TX} be a type-first summary obtained by replacing in Definition 7, \equiv_W by \equiv_X . Based on this counter-example, one can show that the shortcut does not hold for \equiv_{TX} . If the ontology only features subclass triples, the shortcut holds also for \equiv_{TW} and \equiv_{TS} ; this is because any node typed in G^∞ was already typed in G .

Based on Theorem 2, we have also established:

Theorem 7. (BISIMILARITY SHORTCUT) *The shortcut holds for the forward (\equiv_{fb}), backward (\equiv_{bw}), and forward-and-backward (\equiv_{fb}) bisimilarity equivalence relations (recalled in Section 3).*

6.4 Relationships between summaries

On any graph G , it follows from our node equivalence relations that any two nodes that are \equiv_S are also \equiv_W . Based on this, it is easy to show 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, $(G/TS)_{/TW} = G_{/TW}$. 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 13 summarizes the main relationships between G , G^∞ , our summaries and bisimilarity-based ones.

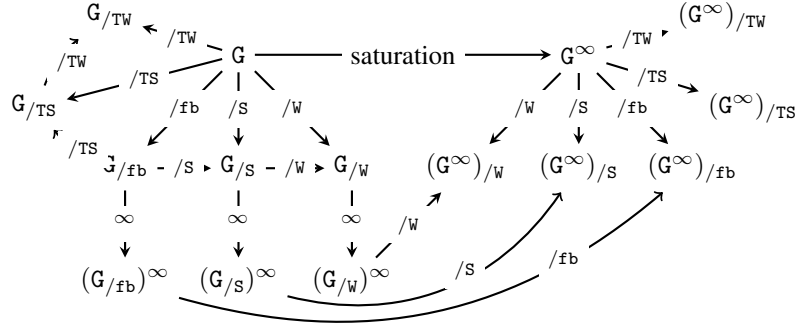


Figure 13: Relations between quotient summaries.

7 Summarization algorithms

We now discuss summarization algorithms which, given as input a graph G , construct G/W , G/S , G/TW and G/TS . The schema S_G is directly copied in the summary thus, below we focus on summarizing D_G (Section 7.1) and T_G (Section 7.2).

7.1 Data graph summarization

We have devised summarization algorithms of two flavors.

Ⓐ **Global** algorithms start by learning the equivalence relation⁶ and creating the summary nodes. Then, a final pass over G computes f_{\equiv} and adds to the summary the edge $f_{\equiv}(n_1) e f_{\equiv}(n_2)$ for each triple $n_1 e n_2$ in G . While this requires several passes over the graph, it connects its nodes directly to their final representatives they will have in $G_{/\equiv}$.

We start with our **global W summarization algorithm** (Table 3). It exploits Proposition 1, which guarantees that any data property occurs only once in the summary. To each distinct data property p encountered in D_G , it associates a summary node (integer) s_p which will be the (unique) source of p in the summary, and similarly a node t_p target of p ; these are initially unknown, and evolve as G is traversed. Further, it uses two maps op and ip which associate to each D_G node n , the set of its outgoing, resp. incoming data properties. These are filled during the first traversal of D_G (step 1.) Steps 2. to 2.5 ensure that for each node n having outgoing properties and possibly incoming ones, s_p for all the outgoing ones are equal, and equal also to t_p for all the incoming ones. This is performed using a function $fuse$ which, given a set of summary nodes, picks one that will replace all of them. In our implementation, summary nodes are integers, and $fuse$ is simply min ; we just need $fuse$ to be distributive over \cup , i.e., $fuse(A, (B \cup C)) = fuse(fuse(A, B), fuse(A, C))$. Step 3. symetrizes this to ensure that the incoming properties of nodes lacking outgoing properties (thus, absent from op) also have the same target. In Step 4., we represent s and o based on the source/target of the property p connecting them. The $fuse$ operations in 2. and 3. have ensured that, while traversing G triples in 4., a same G node n is always represented by the same summary node $f_w(n)$.

Our **global S summarization algorithm** (Table 4) uses two maps sc and tc which store for each data node $n \in D_G$, its source clique $sc(n)$, and its target clique $tc(n)$, and for each data property p , its source clique src_p and target clique trg_p . Further, to each (source clique, target clique) pair encountered until

⁶Recall that $\equiv_w, \equiv_s, \equiv_T$, as well as bisimilarity equivalence, are defined based on the data/type triples of a given graph G , thus when starting to summarize G , we do not know whether any two nodes are equivalent; the full \equiv is known only after inspecting all G triples.

Algorithm global-W(G)

1. For each $s \ p \ o \in G$, add p to $op(s)$ and to $ip(o)$.
2. For each node $n \in op$:
 - 2.1. Let $X \leftarrow fuse\{s_p \mid p \in op(n)\}$.
If X is undefined, let $X \leftarrow nextNode()$;
 - 2.2. Let $Y \leftarrow fuse\{t_p \mid p \in ip(n)\}$.
If Y is undefined, let $Y \leftarrow nextNode()$;
 - 2.3. Let $Z \leftarrow fuse(X, Y)$;
 - 2.4. For each $p \in ip(n)$, let $s_p \leftarrow Z$;
 - 2.5. For each $p \in op(n)$, let $t_p \leftarrow Z$;
3. Repeat 2 to 2.5 swapping ip with op and t_p with s_p ;
4. For each $s \ p \ o \in G$: let $f_w(s) \leftarrow s_p$, $f_w(o) \leftarrow t_p$;
Add $f_w(s) \ p \ f_w(o)$ to G_w .

Table 3: Global W summarization algorithm.

Algorithm global-S(G)

1. For each $s \ p \ o \in G$:
 - 1.1. Check if src_p , trg_p , $sc(s)$ and $tc(o)$ are known; those not known are initialized with $\{p\}$;
 - 1.2. If $sc(s) \neq src_p$, fuse them into new clique $src'_p = sc(s) \cup src_p$; similarly, if $tc(o) \neq trg_p$, fuse them into $trg'_p = tc(o) \cup trg_p$.
2. For each $s \ p \ o \in G$:
 - 2.1. $f_s(s) \leftarrow$ the (unique) summary node corresponding to the cliques $(sc(s), tc(s))$; similarly, $f_s(o) \leftarrow$ the node corresponding to $(sc(o), tc(o))$ (create the nodes if needed).
 - 2.2. Add $f_s(s) \ p \ f_s(o)$ to G_s .

Table 4: Global S summarization algorithm.

a certain point during summarization, we store the (unique) corresponding summary node. Steps 1.-1.2. build the source and property cliques present in G and associate them to every subject and object node (in sc and tc), as well as to any data property (in src_p and trg_p). For instance, on the sample graph in Figure 2, these steps build the cliques in Table 2. Steps 2-2.2. represent the nodes and edges of G .

The correctness of algorithms **global-W** and **global-S** follows quite easily from their descriptions and the summary definitions.

(B) Incremental algorithms simultaneously learn the equivalence relation from G and represent G data triples. They are particularly suited for incremental **summary maintenance**: if new triples Δ_G^+ are added to G , it suffices to run the summarization algorithm only on Δ_G^+ , based on G_{\equiv} and its representation function f_{\equiv} , in order to obtain $(G \cup \Delta_G^+)_{/\equiv}$. Incremental algorithms also provide the basic building blocks for incrementally propagating the effect of a deletion from G . However, incremental algorithms are considerably **more complex**, since various decisions (assigning sources/targets to properties in W , source/target cliques in S , node representatives in both) must be *repeatedly revisited* to reflect newly acquired knowledge. We illustrate this on our algorithms and examples below.

Each **incremental summarization** algorithm consists of an **incremental update method**, called for every D_G triple, which adjusts the summary's data structures, so that at any point, the summary reflects exactly the graph edges (triples) visited until then.

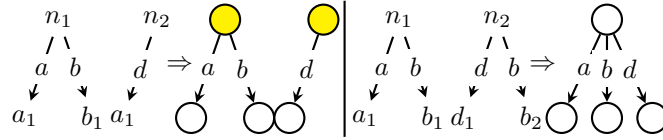
Table 5 outlines incremental W summarization. For example (see the figure below), assume the algorithm traverses the graph G in Figure 2 starting with: $n_1 \ a \ a_1$, then $n_1 \ b \ b_1$, then $n_2 \ d \ d_1$. When we summarize this third triple, we do not know yet that the source of a d triple is also equivalent to n_1 , be-

Algorithm incred-w(s p o)

1. Check if s_p and o_p are known: either both are known (if a triple with property p has already been traversed), or none;
2. Check if $f_w(s)$ and $f_w(o)$ are known; none, one, or both may be, depending on whether s , respectively o have been previously encountered;
3. Fuse s_p with $f_w(s)$ (if one is unknown, assign it the value of the other), and o_p with $f_w(o)$;
4. Update $f_w(s)$ and $f_w(o)$, if needed;
5. Add the edge $f_w(s) p f_w(o)$ to G_w .

Table 5: Incremental W summarization of one triple.

cause no common source of b and d (e.g., n_2 or n_4) has been seen so far. Thus, n_2 is found not equivalent to any node visited so far, and represented separately from n_1 . Now assume the fourth triple traversed is $n_2 b b_2$: at this point, we know that a, b and d are in the same source clique, thus $n_1 \equiv_w n_2$, and their representatives (highlighted in yellow below) must be **fused** in the summary (Step 3.) More generally, it can be shown that \equiv_w **only grows** as more triples are visited, in other words: if in a subset G' of G 's triples, two nodes n_1, n_2 are weakly equivalent, then this holds in any G'' with $G' \subseteq G'' \subseteq G$.



Summary node **fusion** dominates the algorithm's **complexity**. Let N_1, N_2 be two sets of G nodes, represented at a certain point by the distinct summary nodes m_1, m_2 . When fusing the latter into a single m , we must also record that all the nodes in $N_1 \cup N_2$ are now represented by m . A naïve implementation leads to $O(N^2)$ complexity, where N is the number of nodes in D_G , since each new node may lead to a fusion whose cost is $O(N)$; in the worst case N could be proportional to $|G|$, the number of triples in G , leading to an overall complexity of $O(|G|^2)$ for the incremental weak summarization.

Instead, we rely on a **Union-Find** (aka Disjoint Sets) data structure, with the *path compression* and *union by size*⁷ optimizations, which guarantee an overall **quasi-linear worst-case complexity** to our incremental weak summarization algorithm. The exact complexity is $O(N\alpha(N))$ where $\alpha(N)$, the inverse Ackermann's function, is smaller than 5 for any machine-representable input N . Assimilating this to **linear-time**, the algorithm's complexity class is in $O(|G|)$, which is also **optimal**, as summarization cannot do less than fully traversing G .

Table 6 outlines the incremental update of the S summary due to the traversal of the triple $s p o$. Conceptually, the algorithm is symmetric for the source (s) and target (o) of the edge, we only discuss the source side below. Steps 1. and 2. start by determining the source clique of s , based on its previously known source clique (if any) and the previously known source clique of p (if any); after step 2., s 's source (and target) clique reflecting also the newly seen triple $s p o$ are completely known. Determining them may have involved fusing some previously separate cliques. For instance, on the graph in Figure 2, assume we first traverse the two a triples, then we traverse $n_2 b b_2$; so far we have the source cliques $\{a\}$, $\{b\}$ and \emptyset . If the next traversed triple is $n_2 a a_2$, we fuse the source cliques (step 3.1) $\{a\}$ and $\{b\}$ into $\{a, b\}$. This requires fusing the summary node whose (source, target) cliques were $(\{a\}, \emptyset)$ with the one which had $(\{b\}, \emptyset)$ (Step 3.2).

The last intricacy of incremental strong summarization is due to the fact that unlike \equiv_w, \equiv_s **may grow and shrink** during incremental, strong summarization. For instance, assume incremental strong

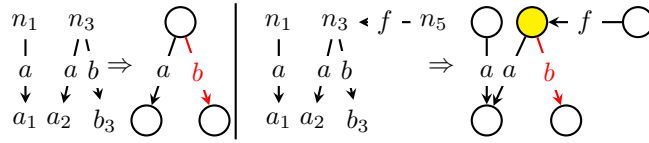
⁷https://en.wikipedia.org/wiki/Disjoint-set_data_structure

Algorithm incred-S(s p o)

1. Check if we already know a source clique src_p (resp. target clique trg_p). Either both are known (if a p triple has already been traversed), or none. Those not known are initialized with $\{p\}$;
2. Check if $sc(s)$ (resp. $tc(o)$) are known; those unknown are initialized with $\{p\}$;
3. If $sc(s) \neq src_p$, fuse them into new clique $src'_p = sc(s) \cup src_p$, using Union-Find; similarly, if $tc(o) \neq trg_p$, fuse them into $trg'_p = tc(o) \cup trg_p$, and:
 - 3.1 Replace $sc(s)$ and src_p with src'_p throughout the summary (respectively, replace $tc(o)$ and trg_p with trg'_p);
 - 3.2 The above may entail summary node fusions; in this case, update f_s (use Union-Find) and the summary edges to reflect it;
4. If before seeing s p o s had been already represented and it had an empty source clique, then s needs to split, i.e., be represented separately from the nodes to which it was \equiv_s previously; call **split-source**(s). (Symetric discussion for o, call **split-target**(o)).
5. Update $f_s(s)$ and $f_s(o)$, if needed;
6. Add the edge $f_s(s) p f_s(o)$ to $G_{/s}$.

Table 6: Incremental S summarization of one triple.

summarization of the graph in Figure 2 starts with $n_1 a a_1, n_3 a a_2, n_3 b b_3$ (see the figure below). After these, we know $n_1 \equiv_s n_3$; their source clique is $\{a, b\}$ and their target clique is \emptyset . Assume the next triple traversed is $n_5 f n_3$: at this point, n_3 is *not* \equiv_s to n_1 *any more*, because n_5 's target clique is now $\{f\}$ instead of the empty \emptyset . Thus, n_5 **splits** from n_1 , that is, it needs to be represented by a new summary node (shown in yellow below), distinct from the representative of n_1 .



Further, note that the representative of n_1 and n_3 (at left above) had one b edge (highlighted in red) which was *solely due to n_3 's outgoing b edge*. By definition of a quotient summary (Section 3), that edge *moves* from the old to the new representative of n_3 (the yellow node). If, above at left, n_1 had also had an outgoing edge labeled b , at right, both nodes in the top row would have had an outgoing b edge. It can be shown that *splits only occur in such cases*, i.e., o whose target clique becomes non-empty (respectively, s whose source clique becomes non-empty, and the node was previously represented together with other nodes; if it was represented alone, we just update the respective clique of its representative).

The procedure **split-source**(s) (omitted for space reasons) represents s separately, to reflect it no longer has an empty target clique, and, for each outgoing edge of s: adds a corresponding edge to the new representative of s; and checks if, as a consequence, an edge needs to be removed from its previous representative.

Proposition 6. (ALGORITHM CORRECTNESS) *Applying algorithm **incred-W** (respectively, **incred-S**) successively on each triple of G , in any order, builds $G_{/w}$ (respectively, $G_{/s}$).*

Splitting requires inspecting the data edges attached to the splitting node, in order to add to its new representative the edges it must have (such as the $n_3 b b_3$ above). We make the hypothesis, denoted (\star), that the maximum number of edges incoming/outgoing a data node is small (and basically constant) compared to the size of G ; this was the case in the graphs we have experimented with. To keep splitting cost under control: (i) We store for each summary node m and edge e a counter $m_{\#}, e_{\#}$ of the D_G nodes

and edges it represents. Splitting is only needed when $m_{\#} > 1$. (ii) Summary node m loses an outgoing (resp. incoming) edge e labeled p when s (resp. o) splits, if and only if the number of outgoing s edges (resp. incoming o edges) labeled p equals $e_{\#}$. At left in the figure above, $e_{\#}$ was 1, thus the b edge is removed from the old representative of n_3 .

Under the (★) hypothesis, using the data structures (including Union-Find) described above, **the complexity of incremental strong summarization is amortised constant per added triple**.

All our algorithms require $O(|G|)$ space to store the summary edges, the representation function, and their other data structures.

7.2 Typed graph summarization

We now explain how to extend our incremental D_G summarization algorithms to type triples.

To extend W , respectively, S summarization to type triples in “data-then-type” fashion (Section 5.1), we run W , resp. S summarization *first, over D_G triples only*, as described in Section 7.1. This assigns their (final) representatives to all nodes of D_G . Then, for each s type C triple, we simply add to the summary the edge $f_w(s)$ type C (resp. $f_s(s)$ type C); recall from Section 5 that any class node C is represented by itself.

For “type-then-data” summarization (Section 5.2), that is, for TW and TS , we *first traverse T_G triples only*, compute all the class sets, and assign to each typed data node a representative based on its class set. Then, we run a *type-aware variant* of a W (resp. S) algorithm, either global or incremental. The changes introduced in the type-aware variant are: (i) In TW summarization, a data property p may lack an untyped source, if p only occurs on typed nodes; consider for instance the graph whose only triples are n_1 type C , $n_1 e a_1$. Similarly, in TS summarization, a property (e.g., e above) may have a target clique, but lack a source clique, since it does not have an untyped source. (ii) Summarizing the data triple $s p o$ does not fuse nor split the representative of s (resp. o) if s (resp. *object*) is typed; instead, representatives of typed nodes (computed during the *type* triples traversal) are left untouched.

Proposition 7. (ALGORITHM CORRECTNESS) *Applying global- W (respectively global- S) on G , or applying increm- W (respectively, increm- S) on each triple of G , extended as described above for data-then-type or type-then-data summarization leads, respectively $G_{/W}$, $G_{/S}$, $G_{/TW}$ and $G_{/TS}$.*

These algorithms need to work with S_G , D_G and T_G *separately*. Fortunately, most popular RDF store allow such direct access. The space needed to also represent T_G triples remains linear in $|G|$.

8 Experimental study

Algorithms and settings. We have implemented our algorithms as well as fw , bw , fb and $1fb$ summarization in Java 1.8. We used the recent [25] algorithm for fw , bw and fb and devised our simple algorithm for $1fb$. We deployed them using PostgreSQL v9.6 to store RDF triples in an integer-encoded triple table, indexed by s , p and o ; the server had 30 GB of shared buffers and 640 MB working memory. The JVM had 90 GB of RAM. We used a Linux server with an Intel Xeon CPU E5-2640 v4 @2.40GHz and 124 GB RAM.

Datasets. We have experimented with numerous real and synthetic datasets, from 10^5 to $1.4 \cdot 10^8$ triples; the largest file is 36.5 GB on disk. Table 7 shows for each graph its number of triples $|G|$, the number of triples in the saturated graph $|G^\infty|$, the number of schema triples $|S_G|$, and the number of distinct data properties $\#p$ and classes $\#C$.

Compression factors. Figure 14 shows the ratio $|G|/|G_{\equiv}|$, called the *compression factor* cf_{\equiv} for our graphs. To our summaries we added $1fb$, fw , bw and fb ; we plot fw and fb , as bw was somewhere in-between. Some fw and bw summarizations ran out of memory or were stopped after more than 2 hours.

| Real datasets | $ G $ | $ G^\infty $ | $ S_G $ | $\#p$ | $\#C$ |
|----------------------|------------|--------------|---------|-------|-------|
| DBLP | 88,153,334 | 88,153,334 | 0 | 26 | 14 |
| DBpedia Person | 7,889,268 | 7,889,268 | 0 | 9 | 1 |
| INSEE Geo | 369,542 | 1,112,803 | 196 | 53 | 144 |
| Springer LOD | 145,136 | 213,017 | 40 | 26 | 16 |
| Nobel Prizes | 80,757 | 109,901 | 35 | 45 | 28 |

| Synth. datasets | $ G $ | $ G^\infty $ | $ S_G $ | $\#p$ | $\#C$ |
|------------------------|-------------|--------------|---------|-------|-------|
| LUBM [16] 1M | 1,227,744 | 1,227,744 | 0 | 17 | 17 |
| LUBM 10M | 11,990,059 | 11,990,059 | 0 | 17 | 21 |
| LUBM 100M | 114,355,295 | 114,355,295 | 0 | 17 | 21 |
| BSBM [2] 1M | 1,000,708 | 1,009,138 | 150 | 38 | 159 |
| BSBM 10M | 10,538,484 | 10,628,484 | 584 | 38 | 593 |
| BSBM 100M | 104,115,556 | 105,315,556 | 2,010 | 38 | 2019 |
| BSBM 138M | 138,742,476 | 140,342,476 | 2,412 | 38 | 2,421 |

Table 7: Datasets used in experiments.

For W and S , cf is close to or above 10^3 , and reaches $3 \cdot 10^6$ in some cases; in all our experiments, cf_w was the highest. In contrast, cf_{fb} rounds to 1 up to 3, since full bisimilarity is rare. Since \equiv_{fw} is a weaker condition and \equiv_{1fb} even weaker, cf_{fw} and especially cf_{1fb} are higher, but still up to 5 times lower than cf_w . *In all our experiments, the weak summary is the most compact; S , TW and TS are close, followed by $1fb$, from 2.9 to 6.9 times larger.* Conversely, full bisimulation achieves little to no compression, while fw (which has the drawback of being asymmetric) compresses less than $1fb$.

Summarization time. The time to build G_w , G_S , G_{TW} and G_{TS} using the global and the incremental algorithms, as well as the time to build G_{1fb} , are plotted as a function of $|G|$ in Figure 15; both axes are in log scale. For each summary type, the summarization time is roughly linear in the size of G , confirming the expectations stated in Section 7; they range from 200 ms to 34.5 minutes (incred-S on the BSBM138M). Incred-W is the fastest overall; it traverses G only once, and it has relatively little processing to do, thus it is faster than global-W which performs several passes, but does not need to replace node representatives. $1fb$ summarization time is close, then we see the global S , TW and TS , and finally incremental S which, as we explained, is quite complex. The fact that increm-W is the cheapest and increm-S the most expensive show that the former may be I/O-bound, while the latter (with the same I/O cost) is CPU (and memory)-bound. Since increm-S is rather expensive per-triple, it is more efficient to first summarize a graph using global-S, and call increm-S only to maintain it later. This is significantly faster: for instance, global-S on BSBM138M takes only 11.85 minutes. Incred-TS is often faster than increm-S because typed nodes do not lead to splits during TS summarization.

Shortcut speed-up. Table 8 shows the time to build $(G^\infty)_{\equiv}$ in two ways: (i) *direct*, i.e., saturate G then summarize, denoted dt_{\equiv} , and (ii) *shortcut*, summarize G , then saturate the summary and summarize again, denoted st_{\equiv} . We define the **shortcut speed-up** x_{\equiv} for $\equiv \in \{W, S\}$ as $(dt_{\equiv} - st_{\equiv})/dt_{\equiv}$ and report it also in the table. The speed-up is highest for G_w (up to almost 98%) and G_S (up to 95%): this is a direct consequence of their good compression. Indeed, dt_{\equiv} includes the time to summarize G^∞ , while st_{\equiv} includes the time to summarize $(G_{\equiv})^\infty$; the smaller this is, the higher x_{\equiv} . The table confirms the practical interest of the shortcut (Section 6) for summarizing the full semantics of a graph G .

Experiment conclusion. Our experiments have shown that our four summaries can be built efficiently, and they reduce graph sizes by factors hundreds up to $3 \cdot 10^6$; they are *two to three orders of magnitude more compact than summaries based on full bisimulation*; G_w is the most compact, and the other three

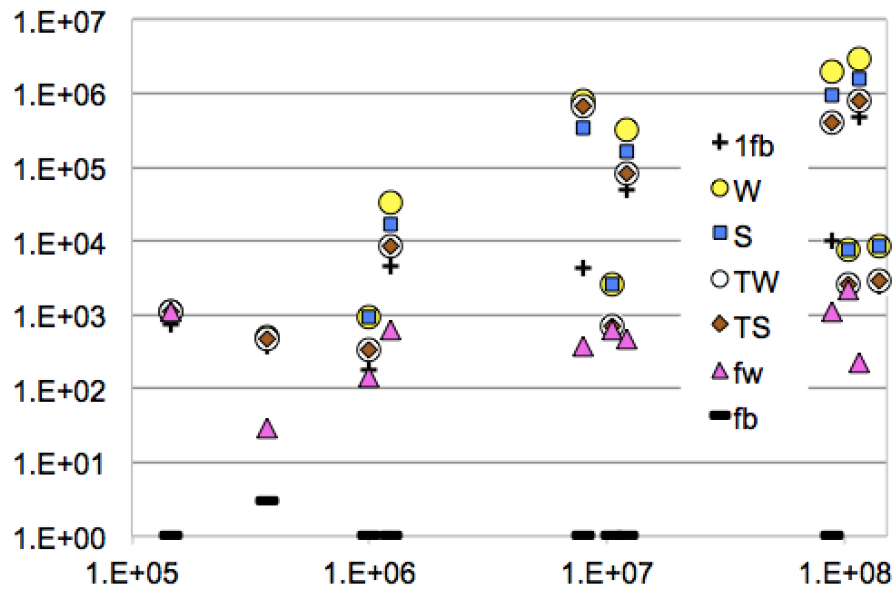
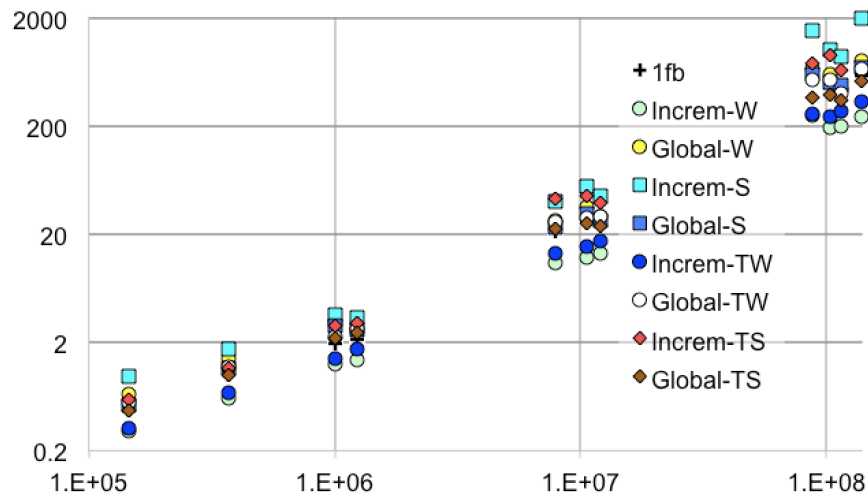


Figure 14: Summarization compression factors.

Figure 15: Summarization time (s) for various graph sizes $|G|$.

summaries we introduced are close (within a factor of 3). Finally, among the summaries which admit a shortcut ($G_{/w}$, $G_{/s}$, $G_{/fb}$ and $G_{/tw}$), the shortcut speed-up is up to 98% for $G_{/w}$, and 95% by $G_{/s}$. All our algorithms scale linearly with $|G|$; increm-W is the fastest, while increm-S is the slowest. Overall, if summary size or building time are a critical concern, we recommend using $G_{/w}$; otherwise, $G_{/s}$ gives finer-granularity information. TW and TS summaries should be used when data types are deemed most important in order to understand the data. However, to summarize G^∞ , only the direct path (saturate, then

| Dataset | dt_w (s) | st_w (s) | x_w (%) | dt_s (s) | st_s (s) | x_s (%) |
|-----------|------------|------------|-----------|------------|------------|-----------|
| INSEE Geo | 35.85 | 0.81 | 97.73 | 38.59 | 1.95 | 94.96 |
| Springer | 3.96 | 0.45 | 88.60 | 4.59 | 1.159 | 74.73 |
| Nobel | 2.13 | 0.41 | 80.56 | 2.60 | 0.75 | 71.09 |
| BSBM1M | 5.45 | 1.48 | 72.85 | 9.35 | 3.82 | 59.20 |
| BSBM10M | 71.63 | 12.67 | 82.32 | 142.46 | 55.64 | 60.94 |
| BSBM100M | 989.00 | 198.00 | 79.98 | 1715.40 | 1030.36 | 39.93 |
| BSBM138M | 1393.69 | 251.93 | 81.92 | 3627.19 | 2049.22 | 43.50 |

Table 8: Shortcut experiments.

summarize) is available for these.

9 Related Work

Graph summarization has a long history and many applications; a recent survey is [28]. It may rely on *graph structure*, *graph values* or *statistics*, e.g., how many distinct values a property has, how many edges are adjacent to a node etc. Summaries may reflect (albeit with some loss of information) the *complete* graph, or they may focus just on *part* of it, e.g., the highest-ranked nodes according to a metric such as PageRank, those having changed recently etc.

In this work, we study *structural quotient summaries*, which are *complete* and *representative* as discussed in Section 3. Our first core contribution is defining property cliques, which lead to the novel notions of weak and strong equivalence, in turn leading to novel quotient summaries for any labeled graph and in particular RDF. Quotient summaries most widely studied in the literature are based on *bisimulation* [29, 9, 21, 12]; they can be built in $O(|G| \cdot \log(N))$, where N is the number of nodes in G . Forward-and-backward bisimulation \equiv_{fb} , symmetric w.r.t. edge directions, is most suited to RDF. However, it is well known, e.g., shown in [25], that heterogeneous graph nodes are very rarely \equiv_{fb} , thus \equiv_{fb} summaries barely compress G (recall Figure 14). \equiv_{1fb} is more permissive, and can be seen as the closest competitor of \equiv_w and \equiv_s ; we have shown \equiv_{1fb} still leads to summaries several times larger than ours. We view our property clique-based summaries as **complementary to** bisimulation-based ones: ours cope better with the heterogeneity of RDF graphs, thus are more suited for visualisation (they were precisely chosen for that in LODAtlas [32]), while bisimulation-based ones lead to larger summaries but allow, e.g., finding “all nodes having properties a and b ” as those represented by a $G_{/1fb}$ node having an a and a b . If one of our summaries is used for this task, a *superset* of the desired nodes is obtained. Bisimulation- and clique-based summaries each have distinct advantages, and can be used side-by-side for different purposes.

The well-known Dataguides [15] are not quotients, as a graph node may be represented by *more than one* node. A Dataguide may be larger than the original graph, and its construction has worst-case exponential time complexity in the size of G . [13] considers *reachability* and *graph pattern* queries on labeled graphs, and builds *answer-preserving summaries* for such queries. However, their query semantics is *not based on graph homomorphisms*, which underlie standard (and SPARQL) query semantics, 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 quotient summaries do preserve.

With a focus farther from our work, [8] introduces an *aggregation* framework for OLAP on labeled graphs, while we focus on representing complete graph structure and semantics. [7] 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 [30, 23]. Quantitative

summarization, where nodes and edges are summarized by their counts, is the focus of [26]. [27] surveys many other quantitative, mining-oriented graph sampling and summarization methods.

Focusing on RDF graphs, [34, 36, 10] study bisimulation-based RDF quotient summaries, providing efficient parallel summarization algorithms [10] and showing they are representative [34]. However, these summaries ignore RDF saturation, and thus its interaction with summarization. Summaries based on clustering [17], user-defined aggregation rules [33], mining [7], and identification of frequent subtrees [38] are not complete and/or require user input. [31] 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; semantics is not studied, either. [3] studies simple methods for summarizing D_G , i.e. the data triples only. Data nodes are distributed in groups, e.g., a node joins a group corresponding to each type (or data property) it has in G ; this entails that some nodes belong in 0 groups, and others may belong to several, unlike the quotient summaries we consider. Implicit triple summarization is not considered.

We had demonstrated [5] and (informally) presented $G_{/w}$ and $G_{/TW}$ in a short “work in progress” paper [4], with procedural definitions (not as quotients). The shortcut was briefly described in a poster [6]. In [18], we propose a type-first summarization technique which exploits subclass hierarchies; beyond the quotient summary framework which it shares, it is orthogonal to the present work.

10 Conclusion

With the goal of obtaining concise yet informative summaries of RDF graphs, we defined property cliques, weak and strong equivalence, provided the first quotient formalization of $G_{/w}$, $G_{/TW}$, and introduced their counterparts $G_{/S}$ and $G_{/TS}$. We have established theoretically (for $G_{/w}$) and experimentally (for all of them) their compactness; they compress their inputs by up to $3 \cdot 10^6$.

Unlike previous equivalence relations mostly focused on *conjunctive* conditions to be satisfied by equivalent nodes, our equivalence relations support a controlled amount of *disjunction*, which makes them much more tolerant to heterogeneity. Thus, while conjunctive summaries (used as structural indexes) can help the query engine evaluate queries, our summaries are complementary; they are meant to help users and application discover the graph structure. This is the use for which they have been integrated in the LODAtlas [32] data exploration and visualization platform; their compactness was also a strong argument for choosing them.

Our second major contribution is the study of the *interplay between saturation and quotient summarization*. We introduced the *shortcut method* for summarizing a semantics-rich RDF graph without saturating it, and formally established a sufficient condition for the shortcut to hold on *any* equivalence relation, based on which we have shown that it holds for \equiv_w , \equiv_S and \equiv_{fb} ; this translates into large performance savings when summarizing G^∞ . Last but not least, we have presented novel, *efficient* $O(|G|)$ algorithms for building our summaries, and for *incrementally maintaining* them in constant per-triple time under the hypothesis (★) (Section 7.1).

Our current work exploits our summaries for summary-based keyword search in RDF graphs; we are also experimenting with parallel summarization algorithms based on Spark.

References

- [1] S. Abiteboul, R. Hull, and V. Vianu. *Foundations of Databases*. Addison-Wesley, 1995.
- [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] Š. Čebirić, F. Goasdoué, and I. Manolescu. A framework for efficient representative summarization of RDF graphs. In *ISWC (poster)*, 2017.
- [7] 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.
- [8] C. Chen, X. Yan, F. Zhu, J. Han, and P. S. Yu. Graph OLAP: towards online analytical processing on graphs. In *ICDM*, 2008.
- [9] Q. Chen, A. Lim, and K. W. Ong. $D(K)$ -index: An adaptive structural summary for graph-structured data. In *SIGMOD*, 2003.
- [10] 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.
- [11] M. P. Consens, R. J. Miller, F. Rizzolo, and A. A. Vaisman. Exploring XML web collections with DescribeX. *TWEB*, 4(3), 2010.
- [12] M. P. Consens, R. J. Miller, F. Rizzolo, and A. A. Vaisman. Exploring XML web collections with DescribeX. *ACM TWeb*, 4(3), 2010.
- [13] W. Fan, J. Li, X. Wang, and Y. Wu. Query preserving graph compression. In *SIGMOD*, 2012.
- [14] F. Goasdoué, I. Manolescu, and A. Roatiş. Efficient query answering against dynamic RDF databases. In *EDBT*, 2013.
- [15] R. Goldman and J. Widom. Dataguides: Enabling query formulation and optimization in semistructured databases. In *VLDB*, 1997.
- [16] Y. Guo, Z. Pan, and J. Heflin. LUBM: A benchmark for OWL knowledge base systems. *J. Web Sem.*, 3(2-3), 2005.
- [17] S. Gurajada, S. Seufert, I. Miliaraki, and M. Theobald. Using graph summarization for join-ahead pruning in a distributed RDF engine. In *SWIM*, 2014.
- [18] P. Guzewicz and I. Manolescu. Quotient RDF Summaries Based on Type Hierarchies. In *DESWeb (Data Engineering meets the Semantic Web) Workshop*, Paris, France, Apr. 2018.
- [19] M. R. Henzinger, T. A. Henzinger, and P. W. Kopke. Computing simulations on finite and infinite graphs. In *FOCS*, 1995.
- [20] R. Kaushik, P. Bohannon, J. F. Naughton, and H. F. Korth. Covering indexes for branching path queries. In *SIGMOD*, 2002.
- [21] R. Kaushik, P. Bohannon, J. F. Naughton, and H. F. Korth. Covering indexes for branching path queries. In *SIGMOD*, 2002.

- [22] R. Kaushik, P. Shenoy, P. Bohannon, and E. Gudes. Exploiting local similarity for indexing paths in graph-structured data. In *ICDE*, 2002.
- [23] K. Khan, W. Nawaz, and Y. Lee. Set-based approximate approach for lossless graph summarization. *Computing*, 97(12):1185–1207, 2015.
- [24] S. Khatchadourian and M. P. Consens. ExpLOD: Summary-based exploration of interlinking and RDF usage in the linked open data cloud. In *ESWC*, 2010.
- [25] S. Khatchadourian and M. P. Consens. Constructing bisimulation summaries on a multi-core graph processing framework. In *GRADES*, 2015.
- [26] K. LeFevre and E. Terzi. GraSS: Graph structure summarization. In *SDM*, 2010.
- [27] S.-D. Lin, M.-Y. Yeh, and C.-T. Li. Sampling and summarization for social networks (tutorial), 2013.
- [28] Y. Liu, T. Safavi, A. Dighe, and D. Koutra. Graph summarization methods and applications: A survey. *ACM Comput. Surv.*, 51(3):62:1–62:34, June 2018.
- [29] T. Milo and D. Suciu. Index structures for path expressions. In *ICDT*, 1999.
- [30] S. Navlakha, R. Rastogi, and N. Shrivastava. Graph summarization with bounded error. In *SIGMOD*, 2008.
- [31] F. Picalausa, Y. Luo, G. H. L. Fletcher, J. Hidders, and S. Vansummeren. A structural approach to indexing triples. In *ESWC*, 2012.
- [32] E. Pietriga, H. Gözükan, C. Appert, M. Destandau, Šejla Čebirić, F. Goasdoué, and I. Manolescu. Browsing linked data catalogs with LODAtlas. In *Int’l. Semantic Web Conference (ISWC), Resources track*, 2018.
- [33] M. Rudolf, M. Paradies, C. Bornhövd, and W. Lehner. SynopSys: large graph analytics in the SAP HANA database through summarization. In *GRADES*, 2013.
- [34] A. Schätzle, A. Neu, G. Lausen, and M. Przyjaciół-Zablocki. Large-scale bisimulation of RDF graphs. In *SWIM*, 2013.
- [35] Y. Tian, R. A. Hankins, and J. M. Patel. Efficient aggregation for graph summarization. In *SIGMOD*. ACM, 2008.
- [36] T. Tran, G. Ladwig, and S. Rudolph. Managing structured and semistructured RDF data using structure indexes. *IEEE TKDE*, 25(9), 2013.
- [37] W3C. Resource description framework. <http://www.w3.org/RDF/>.
- [38] P. Zhao, J. X. Yu, and P. S. Yu. Graph indexing: Tree + delta \geq graph. In *VLDB*, 2007.

For the interested reader, we include here the proofs of all the statements made in the paper. The main proofs are those of our general shortcut Theorems 1 and 2, those of W and S shortcuts (Theorems 3 and 4), and the correctness proofs for our incremental algorithms (Propositions 6 and 7). The other claims proved here are used in these main proofs.

A Proof of Proposition 1

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 corresponds to a set of (disjoint) target cliques STC_1 , which includes TC , and a set of disjoint source cliques SSC_1 . Similarly, n_2 corresponds to a set of (disjoint) target cliques STC_2 , which includes TC , and a set of disjoint source cliques SSC_2 . The presence of TC in STC_1 and STC_2 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

B Proof of Proposition 2

Proof. If two G_S distinct nodes had the same source and the same target clique, they would be strongly equivalent. This cannot be the case in a quotient summary obtained through \equiv_S , since by definition, such a summary has one node for each \equiv_S equivalence class. Thus, any two distinct G_S have distinct source cliques and/or distinct target cliques.

Now, let m be a G_S node, and $S_m = f_S^{-1}(m)$ be the set of all G nodes represented by m . By the definition of a quotient summary, m must be the target (resp. the source) of an edge carrying each of the labels on the edges entering (resp. going out of) any node $n \in S_m$. Thus, m is source of all the properties in the source clique shared by the nodes in S_m , and is target of all the properties in the target clique shared by the nodes in S_m . Thus, m has the source and target clique of any node from S_m ; this concludes our proof. \square

C Proof of Proposition 3

Proof. By definition of an equivalence relation between nodes, every class (resp. property) node is only equivalent to itself, and is represented by itself (i.e., has same label in a graph and in the corresponding summary). Therefore, by definition of a quotient graph using such a node equivalence relation, every edge between class and/or property nodes begets an edge with same label in the graph summary. As schema edges only connects such nodes, schema triples are copied by summarization from a graph to its summary. \square

D Proof of Theorem 1

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

Observe that RDF saturation with RDFS constraints only adds edges between graph nodes, but does not add nodes. Thus, a node n is in G^∞ iff n is in G . Further, by the definition of our quotient-based summaries, n is in G iff $f_\equiv(n)$ is in G/\equiv . Finally, again by the definition of saturation, $f_\equiv(n)$ is in G/\equiv iff $f_\equiv(n)$ is in $(G/\equiv)^\infty$.

Therefore, every G^∞ node n maps the $f_\equiv(n)$ $(G/\equiv)^\infty$ node (*).

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

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

- $n_1 p n_2$ is an edge in G . This holds iff $f_{\equiv}(n_1) p f_{\equiv}(n_2)$ is an edge in $G_{/\equiv}$, by definition of an RDF summary. Finally, if $f_{\equiv}(n_1) p f_{\equiv}(n_2)$ is an edge in $G_{/\equiv}$, then $f_{\equiv}(n_1) p f_{\equiv}(n_2)$ is also an edge in $(G_{/\equiv})^{\infty}$.
- $n_1 p' n_2$ is an edge in G , and p' subproperty p is in S_G^{∞} , thus $n_1 p n_2$ is produced by saturation in G^{∞} . In this case, we show similarly to the preceding item that $f_{\equiv}(n_1) p' f_{\equiv}(n_2)$ is an edge in $(G_{/\equiv})^{\infty}$, hence $f_{\equiv}(n_1) p f_{\equiv}(n_2)$ is also an edge added to $(G_{/\equiv})^{\infty}$ by saturation, since $(G_{/\equiv})^{\infty}$ and G^{∞} have the same (saturated) schema triples (Proposition 3).

If n_1 type c is an edge in G^{∞} , at least one of the following two situations holds:

- n_1 type c is an edge in G . This holds iff $f_{\equiv}(n_1)$ type c is an edge in $G_{/\equiv}$, by definition of an RDF summary (recall that $f_{\equiv}(c) = c$ for classes). Finally, if $f_{\equiv}(n_1)$ type c is an edge in $G_{/\equiv}$, then $f_{\equiv}(n_1)$ type c is also an edge in $(G_{/\equiv})^{\infty}$.
- $n_1 p n_2$ is an edge in G and p domain c (or p range c) is in S_G^{∞} , thus n_1 type c is produced by saturation in G^{∞} . In this case, we show similarly as above that $f_{\equiv}(n_1) p f_{\equiv}(n_2)$ is an edge in $(G_{/\equiv})^{\infty}$, hence $f_{\equiv}(n_1)$ type c is also an edge added to $(G_{/\equiv})^{\infty}$ by saturation, since $(G_{/\equiv})^{\infty}$ and G^{∞} have the same (saturated) schema triples (Proposition 3).

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

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

E Proof of Theorem 2

Proof. We start by introducing some **notations** (see Figure 10). Let f_1 be the representation function from G^{∞} into $(G^{\infty})_{/\equiv}$, and f_2 be the representation function from $(G_{/\equiv})^{\infty}$ into $((G_{/\equiv})^{\infty})_{/\equiv}$.

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

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

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

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

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

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

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

For every edge $n'_1 p n'_2$ in $(G^{\infty})_{/\equiv}$, by definition of an RDF summary, there exists an edge $n_1 p n_2$ in G^{∞} such that $n'_1 p n'_2 = f_1(n_1) p f_1(n_2)$. Figure 16 illustrates the discussion. Further, if $n_1 p n_2$ is in G^{∞} , then $f(n_1) p f(n_2)$ is in $(G_{/\equiv})^{\infty}$ (Theorem 1), hence $f_2(f(n_1)) p f_2(f(n_2))$ is in $((G_{/\equiv})^{\infty})_{/\equiv}$. Therefore,

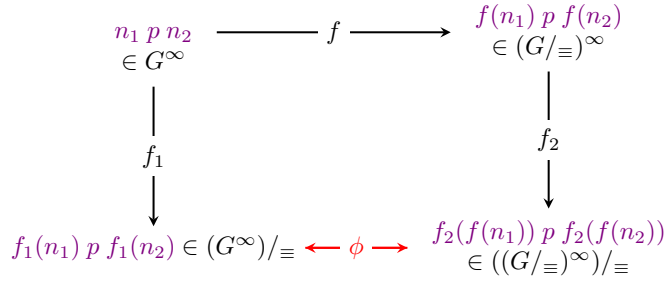


Figure 16: Diagram for the proof of Theorem 2.

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

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

Similarly, for every edge $n'_1 \tau c$ in $(G^\infty)/\equiv$, by definition of an RDF summary, there exists an edge $n_1 \tau c$ in G^∞ such that $n'_1 \tau c = f_1(n_1) \tau c$. Further, if $n_1 \tau c$ is in G^∞ , then $f(n_1) \tau c$ is in $(G/\equiv)^\infty$ (Theorem 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^∞ node, is a bijective function from all $(G^\infty)/\equiv$ nodes to all $((G/\equiv)^\infty)/\equiv$ nodes,
- it follows that $((G/\equiv)^\infty)/\equiv$ contains the image of all $(G^\infty)/\equiv$ $f_1(n_1) \tau c$ triples through φ (*).

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

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

F Proof of Lemma 1

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

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

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

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

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

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

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

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

G Proof of Lemma 2

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

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

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

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

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

H Proof of Proposition 4

Proof. Recall from Lemma 1 that:

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

for some G/W source cliques SC_W^1, \dots, SC_W^m and target cliques TC_W^1, \dots, TC_W^n .

Lemma 2 ensures that the data properties in $(SC_W^1)^+ \cup \dots \cup (SC_W^m)^+$ are related in G^∞ , and those of $(TC_W^1)^+ \cup \dots \cup (TC_W^n)^+$ are related in G^∞ .

Moreover, n_W was created in G/W from a set of weakly-equivalent G nodes having as source clique one among SC_W^1, \dots, SC_W^m and as target clique one among TC_W^1, \dots, TC_W^n . In G^∞ , these nodes connect the data properties of SC_W^∞ with those of TC_W^∞ . \square

I Proof of Theorem 3

Proof. We show that W summaries enjoy the sufficient condition for completeness stated in Theorem 2, i.e., given two nodes n_1, n_2 in G^∞ , f the representation homomorphism corresponding to the weak-equivalence relation \equiv_W , and $f(n_1), f(n_2)$ the images of n_1, n_2 in $(G/W)^\infty$ through f (recall Theorem 1), it holds that: $n_1 \equiv_W n_2$ in G^∞ iff $f(n_1) \equiv_W f(n_2)$ in $(G/\equiv_W)^\infty$.

“Only if”: $n_1 \equiv_W n_2$ in G^∞ iff they are connected by an alternating chain of source and target cliques of G^∞ (as shown in Figure 4); to reuse that figure for the current proof, let us use n_{2k} to denote the n_2 of the current lemma statement. Note that G^∞ only adds triples not nodes, thus all the nodes shown in the figure also exist in G . Now, let us consider the G_W nodes $f(n_1), f(n_2), \dots, f(n_{2k})$ obtained by applying the representation function f on n_1, \dots, n_{2k} . By Theorem 1, f is also a homomorphism from G^∞ to $(G_W)^\infty$, therefore any incoming (outgoing) edge into (from) a node n_j of G^∞ is also incoming (resp. outgoing) into (from) the respective node $f(n_j)$ of $(G_W)^\infty$. As a consequence, we can reproduce

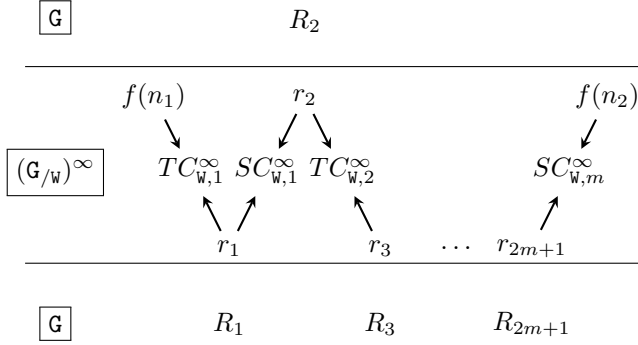


Figure 17: Sketch for sufficient condition for weak summary completeness.

the alternating clique structure into $(G_w)^\infty$, which suffices to make $f(n_1)$ and $f(n_{2k})$ weakly equivalent in $(G_w)^\infty$.

“If”: $f(n_1) \equiv_w f(n_2)$ in $(G_w)^\infty$ iff they are connected by an alternating chain of source and target cliques in $(G_{\equiv_w})^\infty$. Assume w.l.o.g. that the chain is as shown in Figure 17, that is, of the form:

- $f(n_1)$ shares a target clique $TC_{w,1}^\infty$ with r_1
- $r_1(TC_{w,1}^\infty, SC_{w,1}^\infty)$
- $r_2(TC_{w,2}^\infty, SC_{w,1}^\infty), \dots$
- $r_{2m+1}(TC_{w,m-1}^\infty, SC_{w,m}^\infty)$, and $f(n_2)$ has the source clique $SC_{w,m}^\infty$

The alternating chain starts with a target clique and ends with a source clique (of course three other combinations are possible). In the chain, each resource is either:

- $r_{2i+1}(TC_{w,i+1}^\infty, SC_{w,i+1}^\infty)$ or
- $r_{2i+2}(TC_{w,i+2}^\infty, SC_{w,i+1}^\infty)$

for some $0 \leq i < m$. Every r_{2i+1} and r_{2i+2} resource is a node from $(G_w)^\infty$, thus a node from G_w (because saturating G_w does not create nodes). For a given r_j , let R_j be the set of weakly-equivalent G resources from which r_j was created; all resources in R_j are by definition weakly equivalent in G , and this also holds in G^∞ .

By Proposition 4, $TC_{w,1}^\infty$ is also a target clique of G^∞ , and it must be the target clique of n_1 in G^∞ (because of the f homomorphism from G^∞ into $(G_w)^\infty$ ensured by Theorem 1). Similarly, $SC_{w,m}^\infty$ must be a source clique in G^∞ and in particular the source clique of n_2 .

In G^∞ , n_1 shares its target clique $TC_{w,1}^\infty$ with the nodes in R_1 , thus n_1 is weakly-equivalent to any node from R_1 .

Further, by Proposition 4, if the node r_1 has the target clique $TC_{w,1}^\infty$ and the source clique $SC_{w,1}^\infty$ in $(G_w)^\infty$, then the G^∞ node whose target clique is $TC_{w,1}^\infty$ must also have the source clique $SC_{w,1}^\infty$ in G^∞ . (Proposition 4 also ensures that a node in G^∞ has $SC_{w,1}^\infty$ as its source clique.)

If the alternating chain is long enough to comprise r_2 (that is: if the chain does not degenerate in a single node), that corresponds to the set R_2 of G nodes which, in G^∞ , have the source clique $SC_{w,1}^\infty$, therefore they are weakly equivalent to all nodes from R_1 which have the same source clique. Thus, n_1 is weakly equivalent in G^∞ to the nodes from R_1 and R_2 .

The above reasoning can be applied on each edge in the alternating chain, extending weak equivalence from n_1 through all the R_j sets until n_2 . \square

J Proof of Lemma 3

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

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

K Proof of Proposition 5

Proof. Recall from Lemma 1 that:

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

for some S_G source cliques SC_S^1, \dots, SC_S^m and target cliques TC_S^1, \dots, TC_S^n .

Lemma 3 ensures that the data properties in $(SC_S^1)^+ \cup \dots \cup (SC_S^m)^+$ are related in G^∞ , and those of $(TC_S^1)^+ \cup \dots \cup (TC_S^n)^+$ are related in G^∞ .

Moreover, n_S was created in S_G from a set of strongly-equivalent G nodes all sharing a source clique SC_S^i , for $1 \leq i \leq m$, and all sharing a target clique TC_S^j , for some $1 \leq j \leq n$. Thus, in G^∞ , these nodes connect the data properties of SC_S^∞ with those of TC_S^∞ . \square

L Proof of Theorem 4

Proof. “Only if” follows directly from Theorem 1.

To prove “if”, note that $f(n_1) \equiv_S f(n_2)$ in $(S_G)^\infty$ iff they have the same source clique SC_S^∞ and the same target clique TC_S^∞ in $(S_G)^\infty$. By Proposition 5, TC_S^∞ is also a target clique of G^∞ , and it must be the target clique of n_1 and n_2 in G^∞ (because of the f homomorphism from G^∞ into $(S_G)^\infty$ ensured by Theorem 1). Similarly, SC_S^∞ is a source clique in G^∞ , and in particular the source clique of n_1 and n_2 .

Thus, $n_1 \equiv_S n_2$ in G^∞ . \square

M Proof of Theorem 7

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

We show this result using the sufficient condition stated in Theorem 2. That is, $n_1 \equiv_{fw} n_2$ in G^∞ holds iff $f(n_1) \equiv_{fw} f(n_2)$ in $(G/\equiv_{fw})^\infty$ holds.

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

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

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

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

In \mathbb{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 \mathbb{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 $\equiv_{\mathfrak{f}\mathfrak{w}}$ -summary, the set of outgoing properties from $f(n_1)$ to $f(m_1)$ in $\mathbb{G}_{/\equiv_{\mathfrak{f}\mathfrak{w}}}$ 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 $\mathbb{G}_{/\equiv_{\mathfrak{f}\mathfrak{w}}}$ is exactly $\mathcal{P}_{n_2 \rightarrow m_2}$.

Since \mathbb{G} and $\mathbb{G}_{/\equiv_{\mathfrak{f}\mathfrak{w}}}$ have the same schema (Property 3), it follows that in $(\mathbb{G}_{/\equiv_{\mathfrak{f}\mathfrak{w}}})^\infty$, 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 *subproperty* constraints).

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

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

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

In $\mathbb{G}_{/\equiv_{\mathfrak{f}\mathfrak{w}}}$, 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 $\mathbb{G}_{/\equiv_{\mathfrak{f}\mathfrak{w}}}$, 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 $\equiv_{\mathfrak{f}\mathfrak{w}}$ -summary, the set of outgoing properties from n_1 to m_1 in \mathbb{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 \mathbb{G} is exactly $\mathcal{P}_{f(n_2) \rightarrow f(m_2)}$.

Since \mathbb{G} and $\mathbb{G}_{/\equiv_{\mathfrak{f}\mathfrak{w}}}$ have the same schema (Property 3), it follows that in \mathbb{G}^∞ , the set of outgoing properties from n_1 to m_1 , and from n_2 to m_2 , is exactly $\mathcal{P}_{f(n_1), f(n_2) \rightarrow f(m_1), f(m_2)}^\infty$ (data edges can only be added through *subproperty* constraints).

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

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

N Proof of Proposition 6

Proof. All our algorithms (global or incremental) start by identifying the class and property nodes: this is done retrieving all the subjects and objects of S_G triples, and also all the objects of T_G triples. As previously stated, triple stores routinely support such retrieval efficiently. Our algorithms start by representing these special schema nodes exactly by themselves, and copying in the summary all the schema triples from S_G . This exploits Proposition 3 (\mathbb{G} and \mathbb{G}_\equiv have the same schema triples).

Below, we show the correctness of **incremental W and S summarization** on D_G . The proof of Proposition 7 (below) extends this also to T_G triples.

The correctness of **incremental W summarization** on D_G follows from the fact that Algorithm **incremental W** preserves a set of invariants. Let G_k be the first k triples of \mathbb{G} , in the order in which they are traversed by

the algorithm. For any $1 \leq k \leq |G|$, after applying **inrem-w** on k data triples, the following invariants are preserved:

1. The source and target src_p and trg_p of any property p present in these k triples are known.
2. For any summarized triple $s \ p \ o$, we have $f_w(s) = src_p$ and $f_w(o) = trg_p$; further, the summary contains the edge $f_w(s) \ p \ f_w(o)$.

The preservation of these invariants is shown by considering all the cases which may occur for a given summarized triple $s \ p \ o$: the subject s may have *already been seen* (in which case this triple may lead to a fusion), *or not* (in which case we create the new representative of s), and similarly for o . For p there are also two cases (depending on whether we had already encountered it or not, we may create src_p and trg_p , or just fuse them with pre-existing representatives of s and o). There are 8 cases overall. The replacements and fusions detailed in Table 5 guarantee these invariants.

While, for simplicity of presentation, Algorithm **inrem-w** considers the possible fusions due to s and o separately, in reality, given that they may impact the same node(s) (e.g., if $f_w(s) = f_w(o)$), all the replacements are first computed, then reconciled into a *list of summary node substitutions*, applied in all the data structures. For instance, suppose we need to replace summary node 3 with 1 because of a fusion on the subject side, and also summary node 5 with 3 because of a fusion on the object side. In this case, the algorithm will replace 5 and 3 directly with 1. If the replacements were applied sequentially, e.g., first 3 with 1, the second replacement would leave 3 (not 1) instead of 5, which would be an error.

Similarly, the correctness of **incremental S summarization** on D_G follows from the fact that Algorithm **inrem-S** preserves the following invariants after having been called on k successive data triples, with $1 \leq k \leq |G|$:

1. The source and target clique $sc(p)$ and $tc(p)$ of any property p present in these k triples are known, and they contain p .
2. For any summarized triple $s \ p \ o$, we have $f_S(s) = sc(p)$ and $f_S(o) = tc(p)$; further, the summary contains the edge $f_S(s) \ p \ f_S(o)$.
3. For any source clique sc and target clique tc of a node n appearing in the summarized triple, the summary contains exactly one node.
4. For any summary node m , the count $m_\#$ is exactly the cardinality of the set $\{n \in G \mid f_S(n) = m\}$.
5. For any summary edge $m \xrightarrow{p} m'$, the count $e_\#$ is exactly the cardinality of the set $\{n \xrightarrow{p} n' \text{ edge of } G \mid f_S(n) = m \text{ and } f_S(n') = m'\}$.
6. For any (subject, property) combination occurring in the summarized triples, the count $(sp)_\#$ is exactly the number of times this occurred in the triples. Similarly, for any (property, object) combination appearing in the summarized triples, the count $(po)_\#$ is exactly the number of times it appeared.

Like for **inrem-w**, there are eight cases depending on whether s , p and o have been previously seen. Further, in the four cases where s has been seen, we may need to split s 's representative, or not, and similarly for o ; thus, the six cases original cases where at least one of them had been seen lead to 12 cases (to which we add the remaining two, where neither s nor o had been seen), for a total of 14 cases.

Items 4, 5 and 6 are ensured during: the addition of an edge to the summary (this sets $e_\#$ to 1 or increases it); the assignment of representatives to nodes (this sets $m_\#$ to 1 or increments it); the edge repartition during split (this subtracts from one edges $e_\#$ exactly the count that it adds to another new edge); and node replacements (which, when replacing u with v , either carry $u_\#$ into $v_\#$, if v did not exist

in the summary previously, or add $u_{\#}$ to $v_{\#}$ if it did). Together, 4, 5 and 6 ensure the correctness of the **split** algorithm (explained in Section 7.1).

The previous items are ensured by the creation of summary nodes (at most one exists at any time for a given source and target clique), fusing cliques (this guarantees each property is in the right clique, and remove cliques input to the fusion), and replacing / fusing summary nodes, as well as from the correctness of the split procedure. \square

O Proof of Proposition 7

Proof. First, recall that TW and TS summarization start with the type triples, which means all type nodes are detected and represented according to their class sets, before the data triples are summarized. This entails that among the cases which occur for W and S summarization (8, respectively, 14, see discussion in the proof of Proposition 6), those in which the subject, respectively, the object was already represented are further divided in two, depending on whether the subject, respectively, object was a typed node.

This shows that *incremental TW summarization handles a superset of the cases handled by the W one*, and similarly for TS and TS. Thus, **increm-TW**, respectively, **increm-TS** preserve all the invariants of **increm-W**, respectively, **increm-S**⁸, with some additions, which we highlight in italics below:

- Additions of TW summarization w.r.t. W:
 1. The source and target src_p and trg_p of any property p present *with an untyped source, respectively, an untyped target* in the summarized triples are known.
 2. For any summarized triple $s \ p \ o$, we have $f_w(s) = src_p$ *if s is untyped* and $f_w(o) = trg_p$ *if o is untyped*; further, the summary contains the edge $f_w(s) \ p \ f_w(o)$.
- Additions of TS summarization w.r.t. S:
 1. The source and target clique $sc(p)$ and $tc(p)$ of any property p present in these k triples *with an untyped source, respectively, with an untyped target* are known, and they contain p .
 2. For any summarized triple $s \ p \ o$, we have $f_s(s) = sc(p)$ *if s is untyped*, and $f_s(o) = tc(p)$ *if o is untyped*; further, the summary contains the edge $f_s(s) \ p \ f_s(o)$.

Further, they also preserve:

7. The summary contains one node for each set of classes belonging to some resource in the input.
8. For any node n with a non empty class set, $f_{TW}(n)$ (respectively, $f_{TS}(n)$) is the nodes corresponding to the class set of n .

These invariants are ensured by the way in which we collect all class sets during the initial traversal of type triples (common to the TW and TS algorithms). Further, during the TW and TS summarization, as said in Section 7.2, the representatives of typed nodes never fuse, and never split.

The 6 invariants from the above proof of Proposition 6 ensure the correct summarization of D_G triples when s and o are untyped. Together with the two above, they also ensure the correct summarization of triples having a typed source and/or object. \square

⁸Note that in the particular case of triples connecting untyped nodes, the algorithms coincide.



**RESEARCH CENTRE
SOPHIA ANTIPOLIS – MÉDITERRANÉE**

2004 route des Lucioles - BP 93
06902 Sophia Antipolis Cedex

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

ISSN 0249-6399