

# XML-Based RDF Data Management for Efficient Query Processing

Mo Zhou  
Indiana University, USA  
mozhou@cs.indiana.edu

Yuqing Wu  
Indiana University, USA  
yuqwu@cs.indiana.edu

## Abstract

The Semantic Web, which represents a web of knowledge, offers new opportunities to search for knowledge and information. To harvest such search power requires robust and scalable data repositories that can store RDF data and support efficient evaluation of SPARQL queries. Most of the existing RDF storage techniques rely on relation model and relational database technologies for these tasks. They either keep the RDF data as triples, or decompose it into multiple relations. The mis-match between the graph model of the RDF data and the rigid 2D tables of relational model jeopardizes the scalability of such repositories and frequently renders a repository inefficient for some types of data and queries. We propose to decompose RDF graph into a forest of semantically correlated XML trees, store them in an XML repository and rewrite SPARQL queries into XPath/XQuery queries to be evaluated in the XML repository. In this paper, we discuss the basic idea of RDF-to-XML decomposition and the criteria of such decomposition in term of correctness, redundancy and query efficiency, then propose two RDF-to-XML decomposition algorithms based on these criteria. Our experimental evaluation results illustrate that our approach is capable of improving both the storage efficiency and query processing efficiency compared to the existing RDF techniques.

## 1. INTRODUCTION

RDF (Resource Description Framework) [14] is a W3C recommended language for describing linked data of the Semantic Web in the form of triples. RDFS (RDF schema) [7] extends RDF vocabularies to define the structure of underlying RDF data, as well as taxonomic hierarchies of concepts and relations. Both RDF data and schema can be represented as graphs, as shown in Fig. 1. The flexibility, simplicity and expressiveness of the Semantic Web evoke extensive exploitation of RDF/RDFS in many areas such as bioinformatics and social networks.

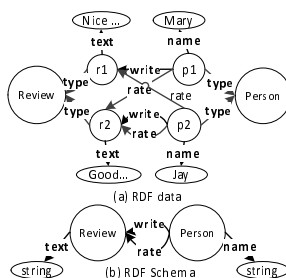


Figure 1: Example RDF data/schema

SPARQL[17] is a W3C recommended RDF query language. A SPARQL query contains a collection of triples with variables called *simple access patterns* which form graph patterns for describing

```
SELECT ?t
WHERE {?p type Person. ?r ?x ?t.
      ?p name ?n. ?p write ?r}
```

query requirements. For example, the SPARQL query on the left retrieves all properties of

reviews written by a person whose name is known.

The needs to develop applications on the Semantic Web and support search in RDF graphs call for RDF repositories to be reliable, robust and efficient in answering SPARQL queries. As in the context of RDB and XML, the selection of storage models is critical to a data repository as it is the dominating factor to determine how to evaluate queries and how the system behaves when it scales up.

### 1.1 Related Works

Most of the existing RDF data repositories [1, 8, 9, 19] rely on relational models for data storage and evaluate SPARQL queries by rewriting them into SQL queries and then executing them in the RDB engine. Among them there are two major directions: (1) keeping the simple triple data model of RDF data, e.g. *triple store* [8, 19]; and (2) decomposing RDF triples into relations, either based on predicates, e.g. *vertical partition* or based on semantics, e.g. *property table* [9].

The *triple store* does not scale well as the evaluation of a complex SPARQL query invokes many self-joins. Various indexing techniques [13, 15, 16, 21] were proposed as remedies, at the cost of huge increase in storage space and decrease in the scalability and update efficiency. The *vertical partition* [1] works well for SPARQL queries when all predicates in the WHERE clause are known. Otherwise, all tables have to be accessed and results unioned. For example, the RDF data in Fig. 1(a) are stored in five tables. All need to be accessed to evaluate the SPARQL query above. The *property table* incurs small number of joins for some queries because a selection in one property table can match multiple simple access patterns. However it suffers storage redundancy and large overhead in query evaluation [1]. For instance, the RDF data in Fig. 1 are stored in two tables, *Person* and *Review*. Both are accessed to evaluate the SPARQL query above because of the predicate variable  $?x$  and variable  $?r$  whose class cannot be determined in the time of query-rewrite.

An alternative approach [3] preserves the graph nature of RDF data by storing RDF graphs in an object-relational database. However, this separates the RDF schema and RDF primary data, which brings difficulties in evaluating queries containing both schema and data instances.

The proposal of serializing RDF graph into XML trees to utilize existing XML technologies [4, 10, 20] focuses on representing all RDF features such as blank node in XML, but pays less or no at-

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

WebDB '10 Indianapolis, IN USA

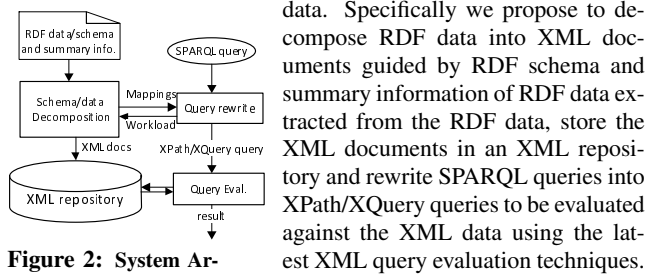
Copyright 2010 ACM 978-1-4503-0186-2/10/06 ...\$10.00.

tention to the efficiency of RDF data storage and query evaluation. It either leads to XML data [20] with large redundancy or flat XML data [10] that cannot take full advantage of XML query evaluation techniques. For example, [10] would serialize the RDF graph in Fig. 1 into a 3-level xml tree and the evaluation of our example query requires eight structural joins and one value join.

## 1.2 Our Proposal

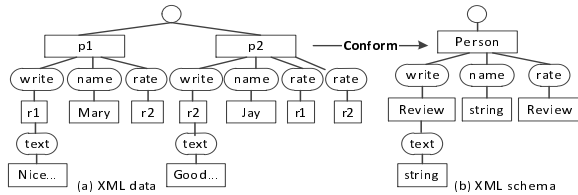
We summarize the desired properties of an RDF storage model to be as follows: (1) *preserving semantics* to facilitate efficient evaluation of RDF queries; (2) *high performance* in evaluating *all* SPARQL queries rather than only some types of SPARQL queries; (3) *high scalability* powered by no or small storage redundancy; and (4) *small overhead* in query rewrite and query optimization.

Semi-structured data model organizes data entries in a tree structure and represents the semantic relationships among them via containment relationships. Tree pattern matching is at the core of the query languages for XML, e.g. XPath[12] and XQuery[11]. We observe the similarity between RDF and XML, in terms of data representation (e.g. using links to represent relationships among data instances) and query (e.g. tree pattern matching in XML and graph pattern matching in RDF), and propose to leverage the sophisticated storage management and query evaluation techniques of XML data repositories, such as [2, 6], to store and query RDF



**Figure 2: System Architecture** that reflects our proposal.

Please consider the example RDF data in Fig. 1. If we decompose it into the XML documents as shown in Fig. 3, the storage requirement is comparable to the relational approaches. Then, the example SPARQL query can be rewritten into an XPath query `//.[name]/write/*//*/*` and be evaluated efficiently due to the high data locality, efficient structural joins and small overhead.



**Figure 3: Example XML data and schema**

RDF data are significantly different from XML data in syntax and data model: RDF data and schema are directed graphs with both nodes and edges labeled, while XML data are trees with only nodes labeled. Although our work, as other RDF storage approaches [1, 3, 8, 9, 19], is syntax independent, the difference between the data models brings substantial challenges to storing and querying RDF data using XML techniques, in transforming graphs into trees, preserving ontology semantics, keeping storage efficiency and mapping graph pattern queries into tree pattern queries.

Our contribution can be summarized as follows:

- We propose an XML-based RDF storage and query evaluation approach which is scalable and capable of supporting various types of SPARQL queries.

- We define the RDF-to-XML decomposition and identify the criteria for measuring the goodness of a decomposition.
- We propose two RDF-to-XML algorithms for decomposing RDF data into XML documents, with the second algorithm taking workload into consideration to further improve the query performance.
- We discuss the query rewrite strategies for rewriting SPARQL queries into XPath/XQuery queries.
- We report on experiments to better understand the storage footprint of our approach, as well as the query performance for various types of SPARQL queries.

## 2. RDF-TO-XML DECOMPOSITION

### 2.1 RDF data and queries

We represent RDF data in the form of a node and edge labeled graph:  $R = (V_R, Ed_R, \lambda_R)$  where  $\lambda_R$  is a labeling function that maps items in  $V_R \cup Ed_R$  into a finite set of labels and literals. We further distinguish RDF schema  $R_s$  from RDF data  $R_d$ , with  $R_s = (V_{R_s}, Ed_{R_s}, \lambda_{R_s})$  representing the class nodes and properties between classes, and  $R_d = (V_{R_d}, Ed_{R_d}, \lambda_{R_d})$  representing the data instances and properties between instances. In particular, in  $R_s$ , we call labels associated with nodes the *class labels* such as *Person* and *Review* in Fig. 1(b), and labels associated with edges the *predicate labels* such as *Write* and *Rate* in Fig. 1(b). Please note that  $V_{R_s} \cup V_{R_d} = V_R$ , but  $Ed_{R_s} \cup Ed_{R_d} \neq Ed_R$  as there exist edges between class nodes and instance nodes in  $R$  indicating the *types* of the instances.

A SPARQL [17] query is frequently represented as graph patterns that are formed by connecting a collection of simple access patterns in the query's WHERE clause. The connection between two simple access patterns with the same variable on their subjects is usually referred to the SS join. Similarly, there are other types of joins, such as SO join, OO join, PP join, etc. Evaluating a SPARQL query can be regarded as mapping the graph patterns against an RDF graph to locate all matches.

### 2.2 Decomposing RDF graphs into XML trees

We propose to decompose RDF graphs into a set of XML trees. We distinguish the decomposition in two steps: (1) the *schema-level decomposition* which maps an RDF schema to a set of XML schemas; and (2) the *data-level decomposition* which maps RDF data to a set of XML documents conforming to the XML schemas obtained from step 1.

We use a directed labeled tree  $(V \cup \{root\}, Ed, \lambda_X)$ , where  $\lambda_X$  is a labeling function that maps nodes in  $V$  to a finite set of labels, to represent an XML schema, and define a mapping from an RDF schema  $R_s$  to a set of XML schemas  $X_s$  as follows:

**DEFINITION 2.1.** *Given an RDF schema  $R_s$  and a set of XML schemas  $X_s = \{X_{s,1}, \dots, X_{s,n}\}$ , the mapping from  $R_s$  to  $X_s$  is a function  $\mathcal{M}_s : R_s \rightarrow X_s$  such that*

- for a node  $v \in V_{R_s}$ ,  $\mathcal{M}_s(v) = \{u | u \in \bigcup_{i=1..n} V_{X_{s,i}} \wedge \lambda_X(u) = \lambda_{R_s}(v)\}$ .
- for an edge  $e = (v_1, v_2) \in Ed_{R_s}$ ,  $\mathcal{M}_s(e) = \{(u_1, u_2) | (u_1, u_2) \in \bigcup_{i=1..n} Ed_{X_{s,i}} \wedge (\lambda_X(u_1) = \lambda_{R_s}(v_1) \wedge \lambda_X(u_2) = \lambda_{R_s}(v_2) \wedge \lambda_X(u) = \lambda_{R_s}(e))\}$ .

Given a mapping  $\mathcal{M}_s$ , we distinguish *class node* and *predicate node* in  $X_s$  based on what RDF schema element (class or predicate) is mapped into it.

Based on the RDF schema decomposition, RDF data  $R_d$  can be transformed into a set of XML documents conforming to  $X_s$  by a

mapping  $\mathcal{M}_d : R_d \rightarrow X_d$ . We say that an XML schema class node *hosts* a data instance in  $R_d$  if the instance is mapped to an XML node conforming to the XML schema class node by  $\mathcal{M}_d$ . For a node  $v \in V_{R_s}$ , if a class node  $u \in \mathcal{M}_s(v)$  can host all instances of the RDF class  $\lambda_R(v)$  in  $R_d$ , we call it a *full node*. Further if  $u$  can preserve the immediate local structure, e.g.  $v$ 's all outgoing edges corresponding to all predicates of the class  $\lambda_R(v)$ , we call  $u$  an *expanded full node* of  $v$ .

Given any class node  $u$  in  $X_s$ , it is desirable to find out whether  $u$  is a full node without computing the mapping  $\mathcal{M}_d$ . In fact, such knowledge can be obtained by analyzing  $R_d$  and summarizing the participation constraint among the classes in  $R_s$ .

**DEFINITION 2.2.** *Given RDF data  $R_d$ , its schema  $R_s$  and two classes  $c_1, c_2$  connected by a predicate  $p$  in  $R_s$ , we say  $c_1$  ( $c_2$ ) is totally participated in a relation  $(c_1, p, c_2)$  (denoted as  $c_1 \vdash_p c_2$  ( $c_2 \vdash_p c_1$ )), if in  $R_d$  any instance of type  $c_1$  ( $c_2$ ) connects to at least one instance of type  $c_2$  ( $c_1$ ) by an edge labeled  $p$ .*

We observe that  $u$  must be a full node if  $u$  is a child of a root in  $X_s$ . In addition, if there are two edges  $(u', e), (e, u) \in \bigcup_{i=1..n} Ed_{X_s, i}$  and  $u'$  is a full node,  $u$  is a full node only if  $\lambda_X(u) \vdash_{\lambda_X(e)} \lambda_X(u')$ .

### 3. DECOMPOSITION CRITERIA

The notions of  $X_s$  and  $X_d$  and the mappings  $\mathcal{M}_s$  and  $\mathcal{M}_d$  are very loosely defined. There are multiple ways to decompose an RDF schema into XML schemas, and therefore multiple ways to decompose RDF data into XML documents. Therefore we propose the prerequisites for the data transformation from RDF to XML to be correctness, low redundancy and high semantic clustering.

While the first criterion must be satisfied, the second and third indeed reflect the trade-off between storage efficiency and query efficiency that is traditional to database system design.

Most RDF schemas have hierarchies where a class has all instances in its subclass and inherits all predicates of its superclasses. However, by analyzing the data, it is always possible to find the finest level of classes whose instances together with relations among them cover the entire dataset but are not overlapped. Therefore, without loss of generality, we only consider such classes and relations among them in the rest of the paper.

#### Correctness

The correctness of the decomposition is determined by the schema decomposition. We say that a set of XML schemas  $X_s$  *structurally covers* a given RDF schema  $R_s$  if any node or edge in  $R_s$  has at least one mapping in  $X_s$ . We say  $X_s$  can *fully cover*  $R_d$  if there exists an  $X_d$  conforming to  $X_s$  such that any node or edge in  $R_d$  has at least one mapping in  $X_d$ .

Instead of taking all possible  $X_s$  fully covering  $R_d$  into consideration, we are interested in a subset that is a straightforward translation of  $R_s$ .

**DEFINITION 3.1.** *Given  $R_s, R_d, X_s, R_d$  and the mappings  $\mathcal{M}_s$  and  $\mathcal{M}_d$ ,  $X_s$  is a full XML schema if for any node  $v \in V_{R_s}$  there exists at least one entry in  $\mathcal{M}_s(v)$  that is an expanded full node.*

#### Redundancy

We summarize two types of redundancy: *Schema Level Redundancy (SLR)* where an RDF schema node is mapped to more than one XML schema node; and *Instance Level Redundancy (ILR)* where an RDF data node is mapped to more than one XML data node that conforms to the same XML schema node.

SLR can be reduced by removing redundant nodes from  $X_s$  under the condition that  $X_s$  still fully covers  $R_s$ . However SLR cannot be fully avoided if  $R_s$  contains circles. It is possible to map any

edge in  $R_s$  to exactly one fragment in  $X_s$  but at least one node in the cycle has to be repeated in  $X_s$ .

The existence of ILR is tied to the cardinality relationship between classes and the predicate connecting them.

**DEFINITION 3.2.** *Given RDF data  $R_d$ , its schema  $R_s$  and two classes  $c_1, c_2$  connected by a predicate  $p$  in  $R_s$ , we say  $(c_1, p, c_2)$  is a 1-1 relation ( $c_1 \leftrightarrow_p c_2$ ), if in  $R_d$  any instance of  $c_1$  is connected to at most one instance of  $c_2$  by  $p$  and vice versa. Similarly, we define 1-n relation ( $c_1 \rightarrow_p c_2$ ) and m-n relation ( $c_1 -_p c_2$ ).*

ILR exists if (1) there exists at least one m-n relation in  $R_s$ ; or (2) two nodes  $c_1$  and  $c_2$  in  $R_s$  with  $c_1 \rightarrow_p c_2$  maps to two class nodes  $c'_1$  and  $c'_2$  in  $X_s$  where  $c'_2$  is an ancestor of  $c'_1$ . ILR in case (1) cannot be eliminated. However ILR caused by case (2) can be eliminated by avoiding such schema mappings.

We say that a decomposition has *minimum redundancy* if no node in  $X_s$  can be removed without losing its full coverage over  $R_d$  and no two XML schema trees in  $X_s$  can be combined into one schema tree without introducing ILR.

#### Query efficiency

We propose to rewrite SPARQL queries into XPath/XQuery queries that are evaluated on XML documents in XML engines. Well-designed RDF-to-XML decomposition strategy can improve the efficiency of query evaluation by increasing data locality, decreasing join evaluation complexity and minimizing the number of joins.

With respect to the data locality we would expect that minimum amount of data needs to be accessed to evaluate a query and the accessed data should be physically stored together. With respect to the joins, value joins can not be avoided when nodes in different XML documents have to be accessed. However, structural joins are much more efficient than value joins [2] and multiple structural joins along a linear pattern can be replaced by a single index access when structural indices are available [6]. Therefore, we strongly favor structural joins over value joins, e.g. favor deeper XML documents than shallower ones, even though it may sacrifice data locality to a certain degree.

## 4. DECOMPOSITION ALGORITHMS

In this section we describe algorithms for decomposing an RDF schema and data into XML schemas and documents that satisfy the criteria in Sec. 3.

### 4.1 Schema Decomposition

We propose two algorithms, greedy decomposition (gR2X) and workload-based decomposition (wR2X), for decomposing an RDF schema into XML schemas. Both algorithms take as inputs an RDF schema  $R_s$  and the summary information reflecting the total participation relations and cardinality relations as discussed above and return XML schemas  $X_s$  as the output. Both are carried out in the following steps:

**Step 1. Initialization.** We create node lists  $V$  and  $F$ .  $V$  contains the class nodes to be assigned to XML schema trees and is initialized to the set of nodes in  $R_s$ . For each class node  $c$  in  $V$ , we create a list  $L_c$  that contains relations  $(c, p, c')$  where  $c \vdash_p c' \wedge c' \rightarrow_p c$ .  $F$  contains the root nodes of the XML schema trees and is initialized to be an empty set.

**Step 2. Tree root selection.** For all nodes  $c$  in  $V$  where  $L_c$  is empty, we create a dummy root node  $r$  in  $F$ , make  $c$  the child of  $r$ , and remove it from  $V$ .

**Step 3. Class node placement.** We maintain a prioritized list of the class nodes in  $V$ . Each time, we pick the node  $c$  at the head of the list, call *chooseRelation*( $L_c$ ) function to pick a relation in  $L_c$  that is to be mapped in the XML schema. If there is no such relation,

we create a dummy root node  $r$  in  $F$ , and make  $c$  the child of  $r$ . Otherwise for each relation  $(c, p, c')$ , we create a predicate node  $u_p$  with label  $p$ , make it the child of  $c'$  and make  $c$  the child of  $u_p$ . In both case we remove  $c$  from  $V$ .

**Step 4. Node expansion.** We expand the class nodes to expanded full nodes so that all predicates that have not been mapped into the XML schema are mapped now and the resultant XML schemas structurally cover  $R_s$ .

Please note that the design of the key of the prioritized list and *chooseRelation* function is the core of the decomposition, as it reflects the strategy on when/where/how a class node is placed.

### Greedy Decomposition (gR2X)

The greedy decomposition uses a greedy approach in *chooseRelation* to construct full XML schemas (Def 3.1) with *minimum redundancy* and deep tree structures. The class nodes with smallest  $L_c$  have higher priority in the prioritized list. Given a class node, the *chooseRelation* function chooses a relation from  $L_c$ , such that the relation  $(c, p, c')$  satisfies the following conditions: 1)  $c'$  is not a descendant of  $c$ ; and 2)  $c'$  is in a tree that is deepest among all the class nodes that have a relation with  $c$  in  $L_c$ . The widest tree is chosen in need of breaking a tie.

### Workload-based Decomposition (wR2X)

It is frequently the case that logical and physical storage models are designed to support efficient evaluation of frequent (sub)queries, which improves the overall throughput. Since every path in the query graph pattern can be mapped to a path in the RDF schema by replacing the nodes along the path by their classes, we use a set of weighted paths in the RDF schema to represent the frequent paths in a workload where weights reflect frequencies. Our workload-based schema decomposition algorithm (wR2X) is to ensure that every frequent path in the workload is represented by a path in the resultant XML schemas.

The wR2X algorithm utilizes the general framework discussed above and takes one extra parameter, a set of weighted paths  $L_e$ , as input. For the prioritized list of class nodes, the key of a class node is the sum of the weights associated with the paths in  $L_e$  containing the node. The *chooseRelation* function chooses relations that are heavily covered by the paths in  $L_e$  to be placed first. More specifically, the *chooseRelation* function chooses a relation from  $L_c$ , so that the relation  $(c, p, c')$  satisfies the following conditions: 1)  $c'$  is not a descendant of  $c$ ; and 2) the sum of the weights associated with paths containing the relation in  $L_e$  is the highest. The deepest tree is chosen in need of breaking a tie.

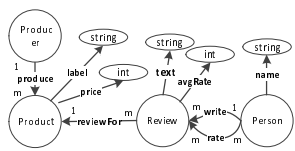


Figure 4: Example RDF schema

After we generate the minimum full XML schemas in step 4, we expand the XML schemas to include all frequent paths in the step as describe below:

**Step 5. Path expansion.** We remove paths in  $L_e$  that are already represented in the XML schemas. We pick the longest path  $lp$  with the highest weight from the remaining paths in  $L_e$ . We then identify the candidate paths in the XML schemas that are rooted at full nodes and contain some fragments of  $lp$ . For each candidate path, we estimate the redundancy that may be introduced if we expand the path to represent  $lp$  and pick the one with the minimum estimated redundancy. We expand the selected path to represent  $lp$  and remove  $lp$  from  $L_e$ . We keep the *remove-pick-expand* process until  $L_e$  is empty.

EXAMPLE 4.1. Consider the example RDF schema shown in Fig. 4. The numbers on the edges are the cardinality information between the corresponding classes. Assume that *Person* is the

only class that does not total participate in any relation, while the other classes total participate in all relations associated with them. Fig. 5 illustrates how the gR2X and wR2X algorithms decompose this RDF schema into corresponding XML schemas<sup>1</sup>. We use the shades of the nodes (light to dark) to indicate in which steps they are added to the resultant schemas.

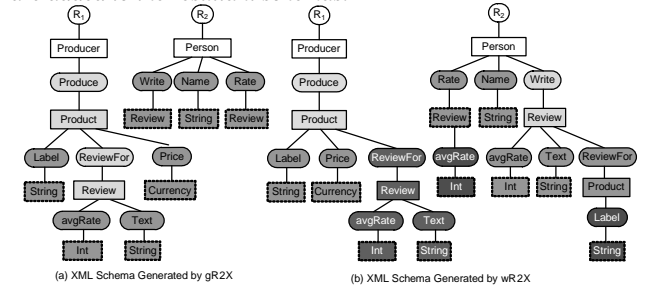


Figure 5: Schema Decomposition

## 4.2 Data Decomposition

Our data decomposition algorithm (R2D) transforms an RDF data  $R_d$  into XML documents conforming to  $X_s$  generated using gR2X or wR2X. Intuitively the transformation process is based on a depth-first traversal of the  $R_d$  graph guided by a depth-first traversal of each XML schema in  $X_s$ . The exact type information is added to an XML node whose type is in a class hierarchy or undeterminable by only referring to the XML schema. Details are omitted due to the space limitation.

EXAMPLE 4.2. Please consider a small fragment of the RDF data shown in Fig. 6(a) whose schema is shown in Fig 4. Fig. 6(b) and (c) show the XML data generated based on the resultant XML schemas in Fig. 5 (a) and (b), respectively.

## 5. QUERY REWRITE

### Query Rewrite Algorithm

We propose to rewrite SPARQL queries into XPath/XQuery queries based on generated XML schemas and the mapping from the RDF schema to these XML schemas. Given a SPARQL query, the rewrite is carried out in three steps.

**Step 1. Class/predicate identification.** We derive the possible classes of each node and predicate variable in the query using the information in the RDF schema, such as the possible predicates of a known class and the range and domain of a known predicate.

**Step 2. Query decomposition.** We decompose the SPARQL query graph into pattern trees that can be matched to subtrees rooted at full nodes in the XML schema. Among many ways of such decomposition we choose the one that satisfies the following criteria: 1) the number of pattern trees are small; 2) the pattern trees feature long paths; and 3) the roots of the pattern trees match to XML schema nodes that are close to the root. We annotate each node in the pattern tree with following information: 1) the role of the node in the query (e.g. variable, optional variable, return variable, URI, or literal), and 2) filtering conditions. In addition, we keep track of the overlapping nodes in these pattern trees as they are the *points of join* to reconstruct the graph pattern from the tree patterns.

**Step 3. XML query construction.** We construct an XPath expression for each pattern tree and assign a variable to it in a *for* clause;

<sup>1</sup>The workload consists of three paths: a. (Producer product Product reviewFor review avgRate Int) with  $w=30$ ; b. (Person rate Review avgRate Int) with  $w=10$ ; and c. (Person write Review ReviewFor Product label String) with  $w=15$

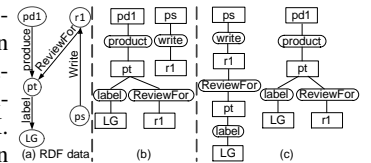


Figure 6: Data Decomposition

we use the *where* to take care of the joins among the pattern trees and the filtering conditions; we then put all return variables into the *return* clause. Hence an XQuery query is constructed.

The rewrite process of the following query is illustrated in Fig. 7. Stratigical construction methods can facilitate query evaluation. For example, XPath2 likely leads to more efficient evaluation plans than XPath1 does by replacing *\*/reviewFor\*/text* to *\*/text*.

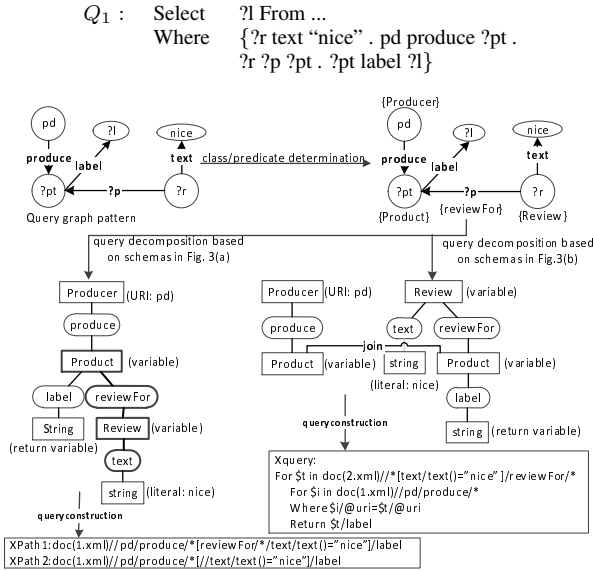


Figure 7: Query rewrite Matched Patterns

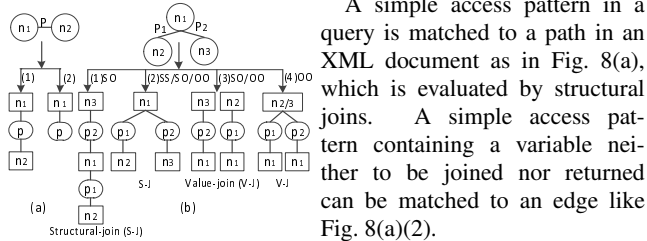


Figure 8: Example Pattern Matching

A simple access pattern in a query is matched to a path in an XML document as in Fig. 8(a), which is evaluated by structural joins. A simple access pattern containing a variable neither to be joined nor returned can be matched to an edge like Fig. 8(a)(2).

Two simple access patterns connected by SS/SO/OO join are matched to either a path or twig(s) in the XML document as in Fig. 8(b). Simple access patterns connected by an SS join are always matched to a twig as Fig. 8(b)(2) and evaluated by structural joins, while those connected with SO/OO joins are matched to path or twig(s) as in Fig. 8(b) and may get value-join involved in the evaluation.

## 6. EXPERIMENTAL EVALUATION

**Experimental Setup** We conducted extensive experiments to compare the XML-based approach we proposed, with the well-know existing approaches, including triple store (TS) [13], vertical partition (VP) [1] and property class table (PT) [9]. We implemented (g/w)R2X, e.g. gR2X and wR2X, on MonetDB/XQuery v4.20.0 and the RDB-based approaches on MonetDB server v4.34.2. We rewrote SPARQL queries into SQL queries for TS, VP and PT following the examples in the Berlin benchmark[5] and the discussion in [18]. We generated XPath/XQuery queries based on the approach in Sec. 5. We timed the hot run [19] of each query. The experiments were carried out on a desktop PC running Ubuntu v8.10

linux (32-bit) with Intel Pentium4 2.80GHz CPU, 3GB memory, and SAMSUNG SP2004C hard drive with 8MB Cache.

**Dataset and Queries** We used the Berlin Benchmark to generate RDF datasets in Turtle format that used compact *in-scope Base URI*. For each dataset, We used *gR2X* and *R2D* algorithms to generate a set of XML documents based on the greedy approach and used *wR2X* and *R2D* with workload that consists of four frequent paths of lengths between 1 and 5 to generate another set of XML documents. In triple store, we sorted triples by PSO and created indices on all permutations of S,P,O; in vertical partition, we sorted tuples by SO and created an index on OS in each table. We report results on seven datasets with 100K, 400K, 1M, 5M, 10M, 15M and 20M<sup>2</sup> triples.

We classified queries in two sets, unbound-variable (UV) queries and extensive-join (EJ) queries. The queries in the UV query set have simple access patterns with variables on nodes or predicates. We report results on four UV queries:  $UV_{s_1/s_2}$  that only have subject variables and  $UV_{p_1/p_2}$  that contain predicate variables. The predicate in  $UV_{s_1}$  is associated to only one class while the one in  $UV_{s_2}$  is associated to many.  $UV_{p_1}$  returns the predicate variable while  $UV_{p_2}$  returns the object variable. Their results are all less than 12 triples.

Query	#SS/OO/SO	#DP	#UC	OP	DC	NR
$EJ_{B_1}$	4/0/0	4	0	N	N	187
$EJ_{B_2}$	0/0/3	14	2	Y	N	20
$EJ_{U_1}$	4/0/0	5	0	N	Y	44
$EJ_{U_2}$	2/0/3	6	1	N	N	134
$EJ_{U_3}$	3/1/1	5	0	N	Y	13875

Table 1: EJ Query Characteristics

The EJ queries feature various number and types of joins. We included various benchmark queries as well as queries we generated to reflect some special features not available in the benchmark queries. We identify the following characteristics of SPARQL queries as factors that have significant impact on the query evaluation: (1) the number of SS/OO/SO joins; (2) the number of distinct properties in the query (#DP); (3) the number of subjects whose classes are neither explicitly indicated by *RDF:type* nor determined by predicates directly associated with them (UC); (4) whether the query contains OPTIONAL operators (OP); (5) whether the query contains variables neither to be joined nor returned (DC); and (6) the cardinality of results (RN). We report results on two benchmark queries ( $EJ_{B_1}$ ,  $EJ_{B_2}$ ) and three queries ( $EJ_{U_1}$ ,  $EJ_{U_2}$ ,  $EJ_{U_3}$ ) of our own. Their characteristics are summarized in Tab. 6. The values

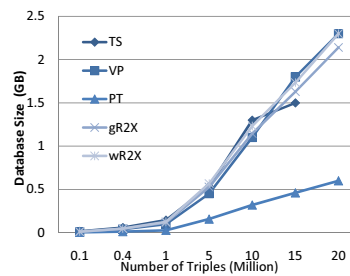


Figure 9: Storage Size

in the NR column reflect the cardinalities of the query results when the queries are evaluated against the dataset with 5M triples.

**Storage Size** We examined the size of physical storage of data with indices and show the results in Fig. 9. PT is most compact since predicates are stored only once in the database catalog. Compared to PT and TS with indices, our (g/w)R2X approaches require slightly less storage and scale in the similar manner.

**Query Evaluation** We conducted experiments on various queries on all data sets we created. Due to the space limitation, we only

<sup>2</sup> MonetDB failed to load 20M dataset into a single table for TS approach. When an XML file exceeded the capacity of the MonetDB/XQuery, we broke it into smaller files. We issued queries against them and unioned the results.

present the query times on the dataset with 5M triples in Fig. 10. We summarize our observations in the following aspects:

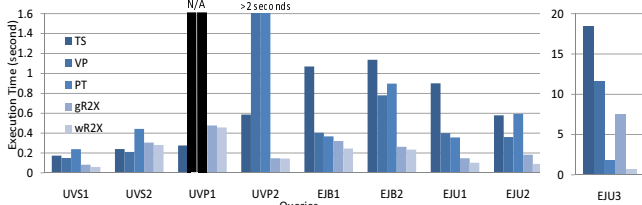


Figure 10: Query Performance Comparison

- (g/x)R2X and TS are capable of evaluating all types of queries purely based on primary data, while PT and VP have to rely on the corresponding database catalog (as in  $UV_{P_2}$ ) or cannot handle such queries without further information (as in  $UV_{P_1}$ )<sup>3</sup>.
- (g/w)R2X delivers good performance in evaluating various queries from simple selection to ones with many joins. TS, VP and PT all favor certain types of queries while having difficulties with others: PT has significant performance retrogression caused by uncertain subject class ( $UV_{S_2}$ ,  $EJ_{B_2}$ ) and predicate variables ( $UV_{P_1}$ ); VP is sensitive to the number of predicates ( $EJ_{B_2}$ ) and predicate variables ( $UV_{P_2}$ ); and TS does not favor complex queries ( $EJ$  queries).
- (g/w)R2X perform very well on path matching queries ( $EJ_{B_2/U_2}$ ) and outperform RDB-based approaches because of the efficient structural joins and high data locality.
- (g/w)R2X access less data than RDB-based approaches, because our query rewrite algorithm restricts the data to be touched in the evaluation process to a small fragment that corresponds to only few classes as discussed in Sec. 5. Only PT has similar feature, but is less precise than our approach in terms of restricting certain matches as the information about the structures is preserved in the mapping and is available to the rewrite algorithm ( $EJ_{U_1}$ ).
- (g/w)R2X can efficiently support queries with OPTIONAL operators because XQuery naturally supports the operator while PT and VP have to involve additional left-outer-join ( $EJ_{B_2}$ ).

Frequently, SPARQL queries feature many simple patterns and mixed types of joins. It is very challenging to evaluate such queries (e.g.  $EJ_{U_3}$ ) against massive data on the Semantic Web. The evaluation of  $EJ_{U_3}$  is very costly for TS and VP due to the loss of the advantage of sort-merge join. Compared to them, PT performs better because of the transformation from SS joins to selections. Structural joins are used in our approaches to evaluate SS and SO joins, while value-join is used here to evaluate OO joins. wR2X performs best among all because of its high data locality. gR2X does not perform well since the deep tree structure chosen by the greedy approach results in lower data locality, but it still outperforms TS and PV significantly.

**Scalability** We studied the scalability of our approaches as well as the RDB-based approaches by comparing the evaluation time of the same query on data of various sizes. We observed that TS scales poorly which is consistent with what was pointed out in the literature. VP scales badly when the query is complex and contains many SO joins. PT scales the best

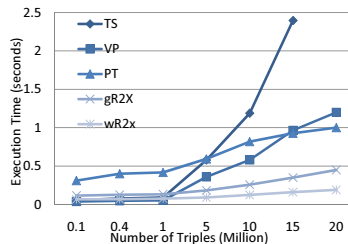


Figure 11: Scalability Study

among the RDB-based approaches but its high parsing and optimization time is too big an overhead for it to be suitable for small datasets. Compared to them, gR2X and wR2X are able to resist the impacts from the complexity of the data. wR2X scales better than gR2X because wR2X preserves data locality in the XML documents even when the dataset scales up. The test result of query  $EJ_{U_2}$  is shown in Fig. 11 to illustrate our observations.

## 7. SUMMARY

To answer the increasing demands on RDF repository, especially high query processing efficiency and scalability, we carefully studied the existing RDF data management systems, identified the preferred properties of an RDF repository and proposed to take advantage of the latest XML data storage and query evaluation techniques by decomposing an RDF graph into XML trees for storage and query evaluation. We identified correctness, low redundancy and high query efficiency as the criteria for a good RDF-to-XML decomposition and proposed two algorithms for decomposing an RDF graph into XML trees. Our experiments demonstrated that compared to other existing RDF storage and query evaluation techniques, our proposed approach requires small storage space, is capable of evaluating SPARQL queries more efficiently and scales better. In addition, our approach is indifferent to the complexity of data and type of queries, making it suitable for supporting the Semantic Web applications in various domains. To further understand, evaluate and improve the approach we proposed, we will compare our approach to the relational approaches at both algorithm and system levels. We will also study how indices and the availability of workload information affect the performance of various systems.

## 8. REFERENCES

- [1] D. J. Abadi, et al. Scalable Semantic Web Data Management Using Vertical Partitioning. In *VLDB*, 2007.
- [2] S. Al-Khalifa, et al. Structural Joins: A Primitive for Efficient XML Query Pattern Matching. In *ICDE*, 2002.
- [3] S. Alexaki, et al. The ICS-FORTH RDFSuite: Managing Voluminous RDF Description Bases. In *SemWeb*, 2001.
- [4] Dave Beckett. RDF/XML syntax specification (revised). W3C Recommendation, February 2004.
- [5] C. Bizer, et al. The Berlin SPARQL Benchmark. In *Int. J. Semantic Web Inf. Syst.*, 5(2): 1-24 (2009).
- [6] S. Brenes, et al. Trie Indexes for Efficient XML Query Evaluation. In *WebDB*, 2008.
- [7] D. Brickley, et al. RDF Vocabulary Description Language 1.0: RDF Schema, Feb 2004.
- [8] J. Broekstra, et al. Sesame: An Architecture for Storing and Querying RDF Data and Schema Information. In *Spinning the Semantic Web*, 2003.
- [9] J. J. Carroll, et al. Jena: implementing the semantic web recommendations. In *WWW*, 2004.
- [10] J. J. Carroll, et al. Rdf triples in xml. In *WWW*, 2004
- [11] D. Chamberlin, et al. XQuery 1.0: An XML query language, May 2003.
- [12] J. Clark, et al. XML path language(XPath), 1999.
- [13] G. H. L. Fletcher, et al. Scalable indexing of RDF graphs for efficient join processing. In *CIKM*, 2009.
- [14] RDF Working Group. Resource Description Framework (RDF), Feb 2004.
- [15] A. Harth, et al. YARS2: A Federated Repository for Querying Graph Structured Data from the Web. In *ISWC/ASWC*, 2007.
- [16] T. Neumann, et al. The RDF-3X engine for scalable management of RDF data. *VLDB J.*, 19(1):91–113, 2010.
- [17] E. Prud'hommeaux, et al. SPARQL Query Language for RDF.
- [18] M. Schmidt, et al. An experimental comparison of rdf data management approaches in a sparql benchmark scenario. In *ISWC*, 2008.
- [19] L. Sidirourgos, et al. Column-store support for RDF data management: not all swans are white. *PVLDB*, 1(2):1553–1563, 2008.
- [20] Norman Walsh. Rdf twig: accessing rdf graphs in xslt. In *Proc. Extreme Markup Languages*, 2003.
- [21] C. Weiss, et al. Hexastore: sextuple indexing for semantic web data management. *PVLDB*, 1(1):1008–1019, 2008.

<sup>3</sup>[18] relies on an additional table that stores all predicates to handle such queries.