

Efficient Evaluation of Nearest Common Ancestor in XML Twig Queries Using Tree-Unaware RDBMS

Klarinda G. Widjanarko, Erwin Leonardi, and Sourav S. Bhowmick

School of Computer Engineering, Nanyang Technological University, Singapore
Singapore-MIT Alliance, Nanyang Technological University, Singapore
{klarinda, lerwin, assourav}@ntu.edu.sg

Abstract. Finding all occurrences of a twig pattern in a database is a core operation in XML query processing. Recent study showed that *tree-aware* relational framework significantly outperform *tree-unaware* approaches in evaluating structural relationships in XML twig queries. In this paper, we present an efficient strategy to evaluate a specific class of structural relationship called *NCA-twiglet* in a tree-unaware relational environment. Informally, *NCA-twiglet* is a subtree in a twig pattern where all nodes have the same nearest common ancestor (the root of *NCA-twiglet*). We focus on *NCA-twiglets* having parent-child relationships. Our scheme is build on top of our SUCXENT++ system. We show that by exploiting the encoding scheme of SUCXENT++ we can reduce useless structural comparisons in order to evaluate *NCA-twiglets*. Through a comprehensive experiment, we show that our approach is not only more scalable but also performs better than a representative tree-unaware approach on all benchmark queries with the highest observed gain factors being 352.

1 Introduction

Finding all occurrences of a twig pattern in a database is a core operation in XML query processing, both in relational implementations of XML databases [3, 6, 7, 8, 12, 13, 14, 19, 20], and in native XML databases [1, 4, 5, 10, 11]. Consequently, in the past few years, many algorithms have been proposed to match twig patterns. These approaches (i) first develop a labeling scheme to capture the structural information of XML documents, and then (ii) perform twig pattern matching based on the labels alone without traversing the original XML documents.

For the first sub-problem of designing appropriate labeling scheme, various methods have been proposed that are based on tree-traversal order [1, 8, 9], region encoding [4, 20], path expressions [10, 14] or prime numbers [17]. By applying these labeling schemes, one can determine the structural relationship between two elements in XML documents from their labels alone. The goal of second sub-problem of matching twig patterns is to devise efficient techniques for structural relationship matching. In general, structural relationship in a twig query may be categorized in two different classes: (a) *NCA-twiglet*, and (b) *path expression*. Given a query twig pattern $Q = (V, E)$, the *nearest common ancestor* (denoted as *NCA*) of two nodes $x \in V, y \in V$ is the common ancestor of x and y whose distance to x (and to y) is smaller than the distance to x of any other common ancestor of x and y . The twig substructure rooted at such *NCA*

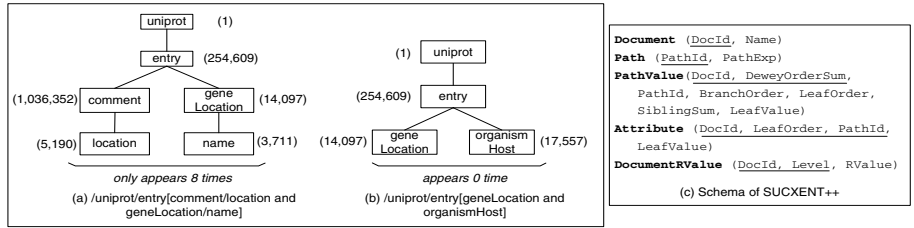


Fig. 1. Example of twig queries and SUCXENT++ schema

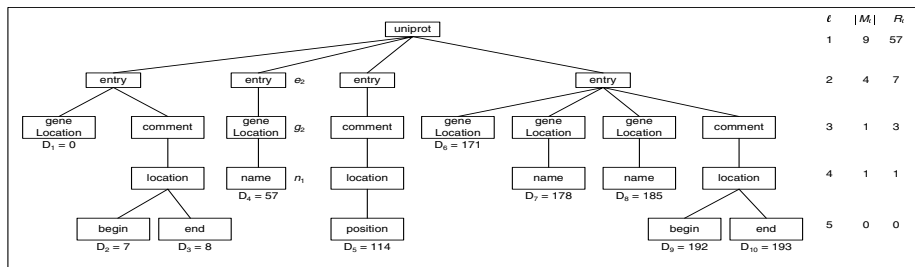


Fig. 2. Example of XML data

node is called *NCA-twiglet*. For example, consider the twig query in Figure 1(a). The twig structure rooted at `entry` node is an example of *NCA-twiglet* as it is the *NCA* of `location` and `name` nodes. On the other hand, *path expression* is a linear structural constraint. For example, `/uniprot/entry` is a *path expression* in Figure 1(a). In this paper, we focus on *efficient evaluation of NCA-twiglets* in a *relational implementation* of XML databases.

In literature, evaluation strategies of twig pattern matching can be broadly classified into the following three types: (a) *binary-structure matching*, (b) *holistic twig pattern matching*, and (c) *string matching*. In the *binary-structure matching* approach, the twig pattern is first decomposed into a set of binary (parent-child and ancestor-descendant) relationships between pairs of nodes. Then, the twig pattern can be matched by matching each of the binary structural relationships against the XML database, and “stitching” together these basic matches [1, 7, 9, 13, 20]. In the *holistic twig pattern matching* approach, the twig query is decomposed into its corresponding path components and each decomposed path component is matched against the XML database. Next, the results of each of the query’s path expressions are joined to form the result to the original twig query [4, 10]. Lastly, approaches like ViST [15] and PRiX [11] are based on *string matching* method and transform both XML data and queries into sequences and answer XML queries through subsequence matching.

A key challenge in *NCA-twiglets* evaluation (as well as twig pattern matching in general) is to develop techniques that can reduce generation of large intermediate results. For instance, the *binary-structure matching* approaches may introduce very large intermediate results. Consider the sample document fragment from UNIPROTKB/ SWISS-PROT and the *NCA-twiglet* in Figures 2 and 1(a), respectively. The path match (e_2, g_2 ,

$n1$) for path `entry/geneLocation/name` does not lead to any final result since there is no `comment/location` path under $e2$. Note that this problem is exacerbated for queries that are *high-selective* but each path in the query is *low-selective*. Note that we use “high-selective” or “very selective” to characterize a twig query with few results and “low-selective” to characterize a query with many results. For example, the query in Figure 1(a) is very selective as it returns only 8 results. However, all the paths are low-selective. The number associated with each node in the queries in Figure 1 represents the number of occurrences of the path from the root node to the specific node in the XML database. Similarly, the query in Figure 1(b) is a high-selective query as it does not return any results although all the paths are low-selective. To solve this problem, the *holistic twig pattern matching* has been developed in order to minimize the intermediate results. In this approach, only those root-to-leaf path matches that will be in the final twig results are enumerated. However, when the twig query contains parent-child (PC) relationships, these solutions may still generate large numbers of useless matches [5]. Hence, in this paper we focus our attention on NCA-twiglets containing *PC relationship* and are components of *high-selective queries having low-selective paths*.

2 Framework and Contributions

The problem of efficiently finding NCAs in a general tree has been studied extensively over the last three decades [2]. Most of these approaches work using some mapping of the tree to a completely balanced tree, thereby exploiting the fact that for completely binary trees the problem is easier. Different algorithms differ by the way they do the mapping. However, these techniques cannot be directly used in the XML context for the following reasons. (i) Although the labels of the nodes used in some of the NCA algorithms can compute the label of NCA in constant time [2], they are not generic enough to efficiently support evaluation of various XPATH axes. Hence, the XML community has resorted to devising novel labeling schemes to support efficient twig matching. (ii) Due to the nature of XML data, the mapping of an XML tree to a completely binary tree may not be an efficient technique for processing different types of XPATH axes. Consequently, the research community has proposed various techniques on native and relational frameworks to evaluate twig queries.

2.1 Relational Approaches for Twig Query Processing and Our Contributions

While a variety of approaches have been proposed in the literature to process twig queries in native XML storage [4, 5, 10, 11], finding ways to evaluate such queries in relational environment has gained significant momentum in recent years. Specifically, there has been a host of work [3, 4, 6, 8, 9, 20] on enabling relational databases to be *tree-aware* by invading the database kernel to implement XML support. On the other side of the spectrum, some completely jettison the approach of internal modification of the RDBMS for twig query processing and resort to alternative *tree-unaware* approach [7, 12, 13] where the database kernel is not modified in order to process XML queries.

While the state-of-the-art tree-aware approaches are certainly innovative and powerful, we have found that these strategies are not directly applicable to relational databases.

The RDBMS systems need to augment their suite of query processing strategies by incorporating special purpose external index systems, algorithms and storage schemes to perform efficient XML query processing. Therefore, the integration of external modules into commercial relational databases could be complex and inefficient. On the other hand, there are considerable benefits in tree-unaware approaches with respect to portability as they do not invade the database kernel. Consequently, they can easily be incorporated in an off-the-shelf RDBMS. However, one of the key stumbling block for the acceptance of tree-unaware approaches has been query performance. In fact, recent results reveal that the tree-aware approaches appear scalable and, in particular, perform orders of magnitude faster than several tree-unaware approaches [3, 8]. In this paper, we explore the challenging problem of *efficient evaluation of NCA-twiglets* in a *tree-unaware* relational framework.

In summary, the main contributions of this paper are as follows. (a) Based on a novel labeling scheme, in Section 3, we present an efficient algorithm for determining nearest common ancestor (NCA) of two elements in an XML document. Our strategy accesses much fewer elements compared to existing state-of-the-art tree-unaware approaches in order to evaluate NCA-twiglets. Importantly, our proposed algorithm is capable of working with any off-the-shelf RDBMS without any internal modification. (b) Through an extensive experimental study in Section 4, we show that our approach significantly outperforms a state-of-the-art tree-unaware scheme (GLOBAL-ORDER [14]) for evaluating benchmark NCA-twiglets.

2.2 Overview of SUCXENT++ Approach

Our approach for NCA-twiglet evaluation is based on the SUCXENT++ system [12]. It is a tree-unaware approach and is designed primarily for query-mostly workloads. Here, we briefly review the storage scheme of SUCXENT++ which we shall be using in our subsequent discussion. The SUCXENT++ schema is shown in Figure 1(c). Document stores the document identifier DocId and the name Name of a given input XML document T . We associate each distinct (root-to-leaf) path appearing in T , namely PathExp, with an identifier PathId and store this information in Path table. For each leaf element n in T , we shall create a tuple in the PathValue table.

SUCXENT++ uses a novel labeling scheme that *does not* require labeling of internal elements in the XML tree. For each leaf element it stores four additional attributes namely LeafOrder, BranchOrder, DeweyOrderSum and SiblingSum. Also, it encodes each level of the XML tree with an attribute called RValue. We now elaborate on the semantics of these attributes. Given two leaf elements n_1 and n_2 , n_1 .LeafOrder $<$ n_2 .LeafOrder iff n_1 precedes n_2 . LeafOrder of the first leaf element in T is 1 and n_2 .LeafOrder = n_1 .LeafOrder+1 iff n_1 is a leaf element immediately preceding n_2 . Given two leaf elements n_1 and n_2 where n_1 .LeafOrder+1 = n_2 .LeafOrder, n_2 .BranchOrder is the level of the NCA of n_1 and n_2 . The data value of n is stored in n .LeafValue.

To discuss DeweyOrderSum, SiblingSum and RValue, we introduce some auxiliary definitions. Consider a sequence of leaf elements C : $\langle n_1, n_2, n_3, \dots, n_r \rangle$ in T . Then, C is a k -consecutive leaf elements of T iff (a) n_i .BranchOrder $\geq k$ for all $i \in [1, r]$; (b) If n_1 .LeafOrder $>$ 1, then n_0 .BranchOrder $<$ k where n_0 .LeafOrder+1 = n_1 .LeafOrder; and (c) If n_r is not the last leaf element in T , then n_{r+1} .BranchOrder $<$ k where

$n_r.\text{LeafOrder}+1 = n_{r+1}.\text{LeafOrder}$. A sequence C is called a *maximal k -consecutive leaf elements* of T , denoted as M_k , if there does not exist a k -consecutive leaf elements C' and $|C| < |C'|$.

Let L_{max} be the largest level of T . The RValue of level ℓ , denoted as R_ℓ , is defined as follows: (i) If $\ell = L_{max} - 1$ then $R_\ell = 1$; (ii) If $0 < \ell < L_{max} - 1$ then $R_\ell = 2R_{\ell+1} \times |M_{\ell+1}| + 1$. For example, consider the XML tree shown in Figure 2. Here $L_{max} = 5$. The values of $|M_1|$, $|M_2|$, $|M_3|$, and $|M_4|$ are 9, 4, 1, and 1, respectively. Then, $R_4 = 1$, $R_3 = 3$, $R_2 = 2 \times 3 \times |M_3| + 1 = 7$, and $R_1 = 2 \times 7 \times |M_2| + 1 = 57$. Note that due to facilitate evaluation of XPATH queries, the RValue attribute in DocumentRValue stores $\frac{R_\ell-1}{2} + 1$ instead of R_ℓ .

DeweyOrderSum is used to encode a element's order information together with its ancestors' order information using a single value. Consider a leaf element n at level ℓ in T . $\text{Ord}(n, k) = i$ iff a is either an ancestor of n or n itself; k is the level of a ; and a is the i -th child of its parent. DeweyOrderSum of n , $n.\text{DeweyOrderSum}$, is defined as $\sum_{j=2}^{\ell} \Phi(j)$ where $\Phi(j) = [\text{Ord}(n, j) - 1] \times R_{j-1}$. For example, consider the rightmost name element in Figure 2 which has a Dewey path "1.4.3.1". DeweyOrderSum of this element is: $n.\text{DeweyOrderSum} = (\text{Ord}(n, 2) - 1) \times R_1 + (\text{Ord}(n, 3) - 1) \times R_2 + (\text{Ord}(n, 4) - 1) \times R_3 = 3 \times 57 + 2 \times 7 + 0 \times 3 = 185$. Note that DeweyOrderSum is not sufficient to compute position-based predicates with QName name tests, e.g., `entry[2]`. Hence, the SiblingSum attribute is introduced to the PathValue table. We do not elaborate further on SiblingSum as it is beyond the scope of the paper.

To evaluate non-leaf elements, we define the *representative leaf element* of a non-leaf element n to be its first descendant leaf element. Note that the BranchOrder attribute records the level of the NCA of two consecutive leaf elements. Let C be the sequence of descendant leaf elements of n and n_1 be the first element in C . We know that the NCA of any two consecutive elements in C is also a descendant of element n . This implies (a) except n_1 , BranchOrder of a element in C is at least the level of element n and (b) the NCA of n_1 and its immediately preceding leaf element is not a descendant of element n . Therefore, BranchOrder of n_1 is always smaller than the level of n . The reader may refer to [12] for details on how these attributes are used to efficiently evaluate ordered XPATH axes.

3 Evaluation of NCA-Twiglets

In this section, we present the evaluation strategy of NCA-twiglets in SUCXENT++. We begin by formally introducing the notion of NCA-twiglet.

3.1 Data Model and NCA-Twiglet

We model XML documents as ordered trees. In our model we ignore comments, processing instructions and namespaces. We also ignore attributes for determining NCA as an attribute is not a child of an element. Queries in XML query languages make use of twig patterns to match relevant portions of data in an XML database. The twig pattern node may be an element tag, a text value or a wildcard "*". We distinguish between query and data nodes by using the term "node" to refer to a query node and the term

“element” to refer to a data element in a document. In this paper, we focus only on parent-child relationships between the nodes in the twig pattern. Recall that existing holistic twig pattern matching approaches achieve optimality for ancestor-descendant relationships but may generate large numbers of useless matches when the twig query contains parent-child relations [5]. We now formally define *NCA-twiglet*.

Definition 1 (NCA-Twiglet). *Given a query twig pattern $Q = \langle V, E \rangle$, a NCA-Twiglet $N = \langle V_n, E_n, \mathfrak{R} \rangle$ in Q , denoted as $N \prec Q$, is a subtree in Q rooted at node $\mathfrak{R} \in V$ such that (a) $V_n \subset V$ is a set of nodes whose nearest common ancestor is \mathfrak{R} , and (b) $E_n \subseteq E$.*

A NCA-twiglet consists of a collection of *rooted path* patterns, where a *rooted path* pattern (RP) is a root-to-leaf path in the NCA-twiglet. The level of the root \mathfrak{R} is called *NCA-level*. For example, the NCA-twiglet in Figure 1(a) consists of the rooted paths `entry/comment/location` and `entry/geneLocation/name`. Note that each of the above RPs has a parent-child relationship between the nodes. The path from $Root(Q)$ to \mathfrak{R} is called the *reachability path* of N . For instance, `/uniprot/entry` is the reachability path.

Given a NCA-twiglet $N \prec Q$ and an XML document D , a match of N in D is identified by a mapping from the nodes in N to the elements in D , such that: (a) the query node predicates are satisfied by the corresponding database elements, wherein wildcard “*” can match any single tag; (b) the parent-child relationship between query nodes are satisfied by the corresponding database elements; and (c) the reachability path of N is satisfied by the database elements. Next, we present our approach to match N in D .

3.2 NCA-Twiglet Matching

Recall that in SUCXENT++ each root-to-leaf path of an XML document is encoded with the attributes `LeafOrder`, `BranchOrder`, `DeweyOrderSum`, and `SiblingSum`. Additionally each level of the XML tree is associated with a *RValue*. Hence, given the NCA-twiglet $N \prec Q$ and document D , our goal is to use these attributes to efficiently determine those root-to-leaf paths that satisfy N . We achieve this by using the following lemma and theorem.

Lemma 1. *Let n_1 and n_2 be two leaf elements in an XML document. If $|n_1.\text{DeweyOrderSum} - n_2.\text{DeweyOrderSum}| < \frac{R_\ell - 1}{2} + 1$ then the level of the nearest common ancestor is greater than ℓ . \square*

Theorem 1. *Let n_1 and n_2 be two leaf elements in an XML document. If $\frac{R_{\ell+1} - 1}{2} + 1 \leq |n_1.\text{DeweyOrderSum} - n_2.\text{DeweyOrderSum}| < \frac{R_\ell - 1}{2} + 1$ then the level of the nearest common ancestor of n_1 and n_2 is $\ell + 1$. \square*

Due to space constraints, we do not present the proof here. The reader may refer to [16] for formal proof. We now illustrate with an example the above lemma and theorem in the context of a twig query. Consider the query in Figure 1(a) and the fragment of the PathValue table in Figure 3 (Step 1). Note that for clarity, we only

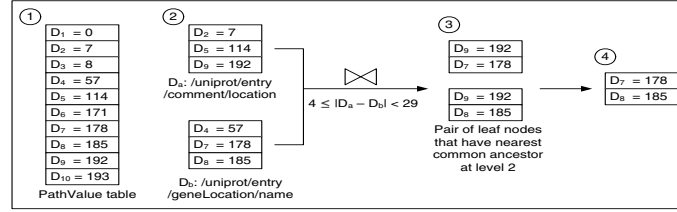


Fig. 3. An example of NCA-twiglet evaluation

<pre> evaluateNCATwiglet (queryTwig) 01 i = 1 02 for every rootedPath in the queryTwig { 03 from_sql.add("PathValue as V_i") 04 where_sql.add("V_i.pathid in rootedPath.getPathId()") 05 where_sql.add("V_i.branchOrder < rootedPath_i.level()") 06 if (i > 1) { 07 where_sql.add("V_i.DeweyOrderSum BETWEEN 08 V_{i-1}.DeweyOrderSum - 09 RValue(rootedPath_i.NCAlevel() - 1) + 1 AND 10 V_{i-1}.DeweyOrderSum + 11 RValue(rootedPath_i.NCAlevel() - 1)") 12 } 13 select_sql.add("DISTINCT V_{i-1}.docId, V_{i-1}.DeweyOrderSum") 14 return select_sql + from_sql + where_sql </pre> <p>(a) evaluateNCATwiglet algorithm</p>	<pre> XPath: /uniprot/entry[comment/location and geneLocation/name] 01 SELECT DISTINCT V2.DocId, V2.DeweyOrderSum 02 FROM PathValue V1, PathValue V2 03 WHERE V1.pathid in (2,3,4) 04 AND V1.branchOrder < 4 05 AND V2.docId = V1.docId 06 AND V2.pathid in (5) 07 AND V2.branchOrder < 4 08 AND V2.DeweyOrderSum BETWEEN 09 V1.DeweyOrderSum - CAST(29 as BIGINT) + 1 AND 10 V1.DeweyOrderSum + CAST(29 as BIGINT) - 1 </pre> <p>(b) An example of Translated SQL query</p>
---	--

Fig. 4. evaluateNCATwiglet algorithm

show the DeweyOrderSums of the root-to-leaf paths in the PathValue table. Let D_a be DeweyOrderSum of the representative leaf elements satisfying `/uniprot/entry/comment/location` (second, fifth, and ninth leaf elements) and D_b be DeweyOrderSum of the representative leaf elements satisfying `/uniprot/entry/geneLocation/name` (fourth, seventh, and eighth leaf elements). This is illustrated in step 2 of Figure 3. From the query we know that D_a and D_b have NCA at level 2 (`/uniprot/entry` level). Hence, based on Theorem 1 we can find pairs of (location,name) elements which have NCA at level 2. D_a and D_b fall on the following range: $(R_2 - 1)/2 + 1 \leq |D_a - D_b| < (R_1 - 1)/2 + 1 \Rightarrow 4 \leq |D_a - D_b| < 29$ which return the (seventh, ninth) and (eighth, ninth) leaf elements pairs (Step 3 of Figure 3). We can easily return the entry subtree by applying Lemma 1 on either one of the elements in the pair (Step 4). Note that since from the XPATH we know that D_a and D_b can not have NCA at level greater than 2, we only need to use Lemma 1 for matching NCA-twiglets. Observe that the above approach can reduce unnecessary comparison as we do not need to find the grandparent of `location` and `name` elements. We can determine the NCA directly by using the DeweyOrderSum and RValue attributes.

3.3 Query Translation Algorithm

Given a query twig (XPATH), the evaluateNCATwiglet procedure (Figure 4(a)) outputs SQL statement. A SQL statement consists of three clauses: `select_sql`, `from_sql` and `where_sql`. We assume that a clause has an `add()` method which encapsulates some simple string manipulations and simple SUCXENT++ joins for constructing valid

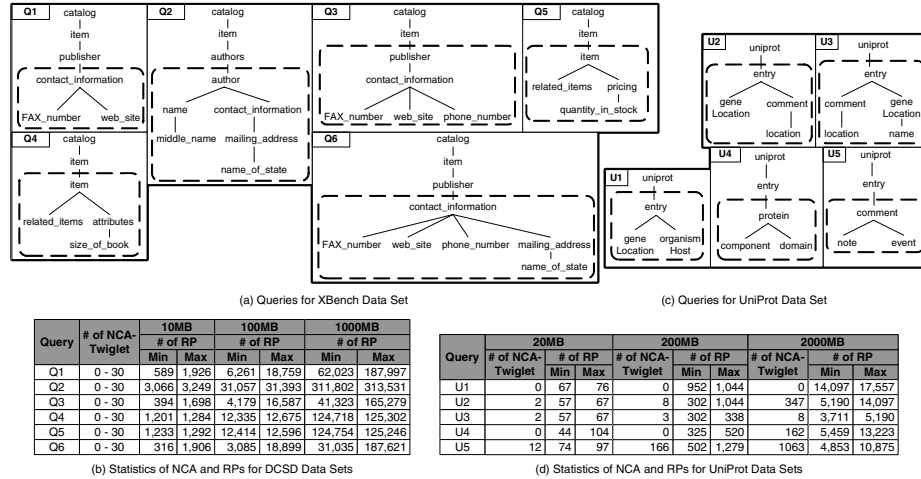


Fig. 5. Query and data sets

SQL statements. In addition to preprocessing PathId, for a single XML document, we also preprocess RValue to reduce the number of joins.

The procedure firstly breaks the query twig into its subsequent rooted path (Line 02). Then for every rooted path, it gets the representative leaf nodes of the rooted path by using PathId and BranchOrder (Lines 04-05). After that, for the second rooted path onwards, it uses Lemma 1 to get the pair of leaf elements that have NCA at the NCA-level (Line 07). After processing the set of rooted paths, we return the DocId and DeweyOrderSum of the rightmost rooted path (Line 11) since only either one of the pairs is needed to construct the whole subtree. Finally, we collect the final SQL statement (Line 12). For example, consider the query in Figure 1(a). The output SQL statement can be seen in Figure 4(b). Lines 03-04 and Lines 06-07 are used to get the representative leaf elements of the respective rooted path. Line 08 is used to get the pair of leaf elements that have NCA at the NCA-level.

4 Performance Study

In this section, we present the performance results of our proposed approach and compare it with a state-of-the-art tree-unaware approach. Since there are several tree-unaware schemes proposed by the community, our selection choice was primarily influenced by the following two criteria. First, the storage scheme of representative approach should not be dependent on the availability of DTD/XML schema. Second, the selected approach must have good query performance for a variety of XPATH axes (ordered as well as unordered) for *query-mostly* workloads. Hence, we chose the GLOBAL-ORDER storage scheme as described in [14]. Prototypes for SUCXENT++ (denoted as SX), and GLOBAL-ORDER (denoted as GO) were implemented with JDK 1.5. The experiments were

conducted on an Intel Pentium 4 3GHz machine running on Windows XP with 1GB of RAM. The RDBMS used was Microsoft SQL Server 2005 Developer Edition.

Data and Query Sets: In our experiments, we used XBench DCSD [18] as synthetic dataset and UNIPROT (downloaded from www.ebi.ac.uk/uniprot/database/download.html) as real dataset. We vary the size of XML documents from 10MB to 1GB for XBench and from 20MB to 2GB for UNIPROT. Recall that we wish to explore twig queries that are high-selective although the paths are low-selective. Hence, we modified XBench dataset so that we can control the number of subtrees (denoted as K) that matches the NCA-twiglet and the number of occurrences of the rooted paths. We set $K \in \{0, 10, 20, 30\}$ for XBench dataset. Note that we did not modify the UNIPROT dataset. Figures 5(a) and 5(c) depict the benchmark queries on XBench and UNIPROT, respectively. We vary the number of rooted paths in the queries from 2 to 4. The number of occurrences of subtrees that satisfies a NCA-twiglet and the minimum and maximum numbers of occurrences of rooted paths in the datasets are shown in Figures 5(b) and 5(d) for XBench and UNIPROT queries, respectively.

Test Methodology: Appropriate indexes were constructed for all approaches through a careful analysis on the benchmark queries. Particularly, for SUCXENT++ we create the following indexes on PathValue table: (a) unique clustered index on PathId and DeweyOrderSum, and (b) non-unique, non-clustered Index on PathId and BranchOrder. Furthermore, since our dataset consists of a single XML document, we removed the DocId column from the tables in SX and GO. Prior to our experiments, we ensure that statistics had been collected. The bufferpool of the RDBMS was cleared before each run. Each query was executed 6 times and the results from the first run were always discarded.

Since GO and SX have different storage approaches, the structure of the returned results are also different. Recall from Section 3.2, the goal of our study is to *identify* subtrees that matches the NCA-twiglet. Hence, we return results in the *select* mode [14]. That is, we do not reconstruct the entire matched subtree. Particularly, for the GO approach, we return the identifier of the root of the subtree (without its descendants) that matches the NCA-twiglet. Whereas for SX, we return the DeweyOrderSum of the root-to-leaf path of the matching subtree. This path must satisfy the rightmost rooted path of the NCA-twiglet. For example, for the query in Figure 1(a), we return the identifiers of the `entry` elements in GO and the DeweyOrderSums of the root-to-leaf paths containing the rightmost rooted path `entry/geneLocation/name` elements in SX. Lastly, for SX we enforce a “left-to-right” join order on the translated SQL query using query hints. The performance benefits of such enforcement is discussed in [12].

NCA-twiglet evaluation times: Our experimental goal is to measure the evaluation time for determining those subtrees that match a NCA-twiglet with a specific reachability path in the twig queries in Figure 5. Figures 6(a) and 6(b) depict the NCA-twiglet evaluation times of SUCXENT++ and GLOBAL-ORDER, respectively. Figure 6(c) depicts the evaluation time for UNIPROT data set. We observe that SX significantly outperforms GO for all queries with the highest observed factor being 352 (Query *U5* on 2GB dataset). Particularly, SX is orders of magnitude faster for high-selective queries. Observe that for XBench dataset, when $K = 0$, SX is up to 332 times faster (Query *Q6* on 1GB dataset)

ID	10MB				100MB				1000MB			
	K=0	K=10	K=20	K=30	K=0	K=10	K=20	K=30	K=0	K=10	K=20	K=30
Q1	25.60	27.00	28.60	27.40	142.80	157.40	151.80	158.20	1,586.80	1,564.80	1,574.20	1,596.60
Q2	61.00	62.40	62.60	69.80	430.40	418.20	475.20	427.00	3,846.80	3,631.20	3,994.80	3,619.60
Q3	68.80	85.80	85.80	84.20	100.00	182.60	185.60	189.00	990.80	1,664.80	1,620.60	1,640.40
Q4	26.00	28.00	27.60	27.40	205.20	168.20	194.40	180.00	1,980.40	1,995.40	1,950.60	1,985.20
Q5	63.80	75.40	62.20	75.60	180.20	186.80	192.80	189.20	1,994.40	1,947.60	1,963.20	1,952.60
Q6	92.00	100.40	112.40	114.00	83.60	257.00	239.20	237.60	617.20	2,161.20	2,159.00	2,159.00

(a) SUCXENT++ (XBench, in msec)

ID	10MB				100MB				1000MB			
	K=0	K=10	K=20	K=30	K=0	K=10	K=20	K=30	K=0	K=10	K=20	K=30
Q1	609.60	603.00	936.40	717.60	8,478.80	8,316.60	8,862.20	6,066.80	114,717.60	83,035.80	80,979.60	84,053.80
Q2	599.00	467.00	443.20	448.80	6,080.00	5,996.40	5,451.40	6,640.00	349,974.00	236,906.20	226,509.20	225,286.20
Q3	698.20	736.00	957.20	675.00	5,510.60	5,565.20	5,458.40	5,464.60	80,954.80	78,541.20	76,998.20	80,571.40
Q4	494.80	474.00	473.80	763.80	3,522.80	4,250.40	4,727.20	4,598.80	107,910.20	108,039.40	71,082.40	111,399.20
Q5	553.60	492.60	570.20	655.60	4,915.40	4,959.40	4,010.60	4,161.60	70,275.20	111,901.60	71,232.80	75,587.40
Q6	762.40	1,205.00	792.60	970.60	7,516.40	6,925.40	7,907.20	7,948.60	204,835.40	249,687.40	259,323.20	238,127.40

(b) Global Order (XBench, in msec)

ID	SUCXENT++			Global Order		
	20MB	200MB	2000MB	20MB	200MB	2000MB
U1	9.00	23.40	163.60	201.00	1,424.60	17,281.60
U2	9.00	21.00	156.80	363.80	2,512.60	23,839.20
U3	5.00	12.20	86.00	266.00	2,675.20	23,705.20
U4	5.40	14.40	123.60	14.00	618.80	6,870.80
U5	5.00	22.20	123.60	438.80	2,800.00	43,502.20

(c) UniProt (in msec)

Fig. 6. Performance results

and on average 56 times faster than GO. This is significant in an environment where users would like to issue exploratory ad hoc queries. In this case, the user would like to know quickly if the query returns any results. If the result set is empty then he/she can further refine his/her query accordingly.

SX is significantly faster than GO because of the following reasons. Firstly, SX uses an efficient strategy based on Theorem 1 to reduce useless comparisons. Furthermore, the number of join operations in GO is more than SX. For example, for $Q6$, GO and SX join six tables and four tables, respectively. Secondly, GO stores every element of an XML document whereas sx stores only the root-to-leaf paths. Consequently, the number of tuples in the Edge table is much more than that in the PathValue table.

5 Related Work

We first compare our proposed approach with existing tree-unaware techniques [7, 12, 13, 14, 19]. Note that we do not compare our work with tree-aware relational schemes [1, 3, 6, 8, 9, 20] as these techniques modify the database internals. Our approach differs from these tree-unaware techniques in the following ways. First, we use a novel and powerful numbering scheme that only encodes the leaf elements and the levels of the XML tree. In contrast, most of the tree-unaware approaches encode both internal and leaf elements. Second, the translated SQL of SUCXENT++ does not suffer from large number of joins. Third, all previous tree-unaware approaches, reported query performance on XML documents with small/medium sizes – smaller than 500 MB. We investigate query performance on large synthetic and real datasets (up to 2GB). This gives more insights on the scalability of the state-of-the-art tree-unaware approaches for twig query processing.

In our previous work [12], we focused on efficiently evaluating ordered path expressions rather than tree-structured queries. In this paper, we investigate how the encoding

scheme in [12] can be used for efficiently processing NCA-twiglet, a specific class of structural relationship in a twig pattern query.

6 Conclusions

The key challenge in XML twig pattern evaluation is to efficiently match the structural relationships of the query nodes against the XML database. In general, structural relationship in a twig query may be categorized in two different classes: path expression and NCA-twiglet. A path expression enforces linear structural constraint whereas NCA-twiglet specifies tree-structured relationship. In this paper, we present an efficient strategy to evaluate NCA-twiglets having parent-child relationship in a tree-unaware relational environment. Our scheme is build on top of SUCXENT++ [12]. We show that by exploiting the encoding scheme of SUCXENT++ we can reduce useless structural comparisons in order to evaluate NCA-twiglets. Our results showed that our proposed approach outperforms GLOBAL-ORDER [14], a representative *tree-unaware* approach for all benchmark queries. Importantly, unlike tree-aware approaches, our scheme does not require invasion of the database kernel to improve query performance and can easily be built on top of any off-the-shelf RDBMS.

References

1. Al-Khalifa, S., Jagadish, H.V., Patel, J.M., et al.: Structural Joins: A Primitive for Efficient XML Query Pattern Matching. In: ICDE (2002)
2. Alstrup, S., Gavaille, C., Kaplan, H., Rauhe, T.: Nearest Common Ancestors: A Survey and a new Distributed Algorithm. In: SPAA (2002)
3. Boncz, P., Grust, T., van Keulen, M., Manegold, S., Rittinger, J., Teubner, J.: MonetDB/XQuery: A Fast XQuery Processor Powered by a Relational Engine. In: SIGMOD (2006)
4. Bruno, N., Koudas, N., Srivastava, D.: Holistic Twig Joins: Optimal XML Pattern Matching. In: SIGMOD (2002)
5. Chien, S., Li, H-G., Tatemura, J., et al.: Twig²Stack: Bottom-up Processing of Generalized-Tree-Pattern Queries over XML Documents. In: VLDB (2006)
6. DeHaan, D., Toman, D., Consens, M.P., Ozsü, M.T.: A Comprehensive XQuery to SQL Translation Using Dynamic Interval Coding. In: SIGMOD (2003)
7. Florescu, D., Kossman, D.: Storing and Querying XML Data using an RDBMS. IEEE Data Engg. Bulletin 22(3) (1999)
8. Grust, T., Teubner, J., Keulen, M.V.: Accelerating XPath Evaluation in Any RDBMS. In: ACM TODS, vol. 29(1) (2004)
9. Li, Q., Moon, B.: Indexing and Querying XML Data for Regular Path Expressions. In: VLDB (2001)
10. Lu, J., Ling, T.W., Chen, T.Y., Chen, T.: From Region Encoding to Extended Dewey: On Efficient Processing of XML Twig Pattern Matching. In: VLDB (2005)
11. Rao, P., Moon, B.: PRIX: Indexing and Querying XML Using Prüfer Sequences. In: ICDE (2004)
12. Seah, B.-S., Widjanarko, K.G., Bhowmick, S.S., Choi, B., Leonardi, E.: Efficient Support for Ordered XPath Processing in Tree-Unaware Commercial Relational Databases. In: DASFAA (2007)

13. Shanmugasundaram, J., Tufte, K., et al.: Relational Databases for Querying XML Documents: Limitations and Opportunities. In: VLDB (1999)
14. Tatarinov, I., Viglas, S., Beyer, K., et al.: Storing and Querying Ordered XML Using a Relational Database System. In: SIGMOD (2002)
15. Wang, H., Park, S., Fan, W., Yu, P.S.: ViST: A Dynamic Index Method for Querying XML Data by Tree Structures. In: SIGMOD (2003)
16. Widjanarko, K.J., Leonardi, E., Bhowmick, S.S.: Efficient Evaluation of Nearest Common Ancestor in XML Twig Queries Using Tree-Unaware RDBMS. Technical Report (2007), Available at <http://www.cais.ntu.edu.sg/~assourav/TechReports/nca-TR.pdf>
17. Wu, X., Lee, M., Hsu, W.: A Prime Number Labeling Scheme for Dynamic Ordered XML Trees. In: ICDE (2004)
18. Yao, B., Tamer Özsu, M., Khandelwal, N.: XBench: Benchmark and Performance Testing of XML DBMSs. In: ICDE (2004)
19. Yoshikawa, M., Amagasa, T., Shimura, T., Uemura, S.: XRel: a path-based approach to storage and retrieval of xml documents using relational databases. ACM TOIT 1(1), 110–141 (2001)
20. Zhang, C., Naughton, J., Dewitt, D., Luo, Q., Lohmann, G.: On Supporting Containment Queries in Relational Database Systems. In: SIGMOD (2001)