

# From Region Encoding To Extended Dewey: On Efficient Processing of XML Twig Pattern Matching

Jiaheng Lu Tok Wang Ling Chee-Yong Chan Ting Chen

Department of Computer Science  
National University of Singapore

{lujiahen,lingtw,chancy,chent}@comp.nus.edu.sg

## Abstract

Finding all the occurrences of a twig pattern in an XML database is a core operation for efficient evaluation of XML queries. A number of algorithms have been proposed to process a twig query based on *region encoding* labeling scheme. While region encoding supports efficient determination of structural relationship between two elements, we observe that the information within a single label is very *limited*. In this paper, we propose a new labeling scheme, called *extended Dewey*. This is a *powerful* labeling scheme, since from the label of an element alone, we can derive all the elements names along the path from the root to the element. Based on *extended Dewey*, we design a novel holistic twig join algorithm, called TJFast. Unlike all previous algorithms based on region encoding, to answer a twig query, TJFast only needs to access the labels of the *leaf* query nodes. Through this, not only do we reduce disk access, but we also support the efficient evaluation of queries with wildcards in branching nodes, which is very difficult to be answered by algorithms based on region encoding. Finally, we report our experimental results to show that our algorithms are superior to previous approaches in terms of *the number of elements scanned, the size of intermediate results and query performance*.

---

*Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the VLDB copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Very Large Data Base Endowment. To copy otherwise, or to republish, requires a fee and/or special permission from the Endowment.*

**Proceedings of the 31st VLDB Conference,  
Trondheim, Norway, 2005**

## 1 Introduction

With the increasing popularity of XML for data representation, there is a lot of interest in query processing over data that conforms to a *tree-structured* data model. Queries on XML data are commonly expressed in the form of tree patterns (or twig patterns), which represent a very useful subset of XPath and XQuery. Efficiently finding all twig pattern matches in an XML database is a major concern of XML query processing. In the past few years, many algorithms ([3],[6],[11],[10]) have been proposed to match such 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 a proper labeling scheme, various methods have been proposed that are based on *tree-traversal* order (e.g. extended preorder [12]), textual positions of the *start* and *end* tags (e.g. region encoding [3]), path expressions (e.g. Dewey ID [22], PID [2]) or prime numbers (e.g. [25]). By applying these labeling schemes, one can determine the relationship (e.g. ancestor-descendant) between two elements in XML documents from their labels alone. Although existing labeling schemes preserve the positional information within the hierarchy of an XML document, we observe that the information contained by a single label is very *limited*. As an illustration, let us consider the most popular *region encoding* scheme, where each label consists of a 3-tuple (*start, end, level*). Given the labels of two elements, one can determine how the elements are structurally related (i.e. ancestor-descendant, parent-child relationships). However, the information derived from a single label is very limited. For instance, the label does not provide any information about the name (i.e. type) of any element.

In this paper, motivated by the existing *Dewey ID* [22], we propose a new *powerful* labeling scheme, called

*extended Dewey ID* (for short, *extended Dewey*). The unique feature of this scheme is that, from the label of an element alone, we can *derive the names of all elements in the path from the root to this element*. For example, Figure 1 shows an XML document with *extended Dewey* labels. Given the label “0.5.1.1” of element *text* alone, we can derive that the path from the *root* to *text* is “/bib/book/chapter/section/text”. An immediate benefit of this feature is that, to evaluate a twig pattern, we *only need to access the labels of elements that satisfy the leaf node predicates in the query*. Further, this feature enables us to easily match a path pattern by string matching. Take element “0.5.1.1” as an example again. Since we see that its path is “/bib/book/chapter/section/text”, it is quite straightforward to determine whether this path matches a path query (e.g. “//section/text”). As a result, the *extended Dewey* labeling scheme provides us an *extraordinary* chance to develop a new efficient algorithm to match twig patterns.

For the second sub-problem of performing structural joins efficiently, several algorithms have been developed to process twig queries. In particular, Bruno et al. [3] proposed the holistic twig matching algorithms PathStack/TwigStack. For evaluating queries with only *ancestor-descendant*(A-D) edges, TwigStack guarantees that each intermediate path solution contributes to final answers. Lu et al.([13]) proposed TwigStackList to efficiently handle twig queries with *parent-child*(P-C) relationships.

Wildcard steps in XPath are commonly used when element names are unknown or do not matter([5]). Previous holistic twig matching algorithms are inefficient for queries with wildcards in branching nodes. For example, consider the XPath query: //a/\*[b]/c. By knowing only the region encodings of *a*, *b* and *c*, we cannot answer this query.<sup>1</sup> How can we answer such queries efficiently?

In this paper, we propose a novel holistic twig join algorithm, called TJFast(i.e. a Fast Twig Join algorithm) based on *extended Dewey* labeling scheme. To match a twig pattern, our algorithm only scans elements for query *leaf* nodes. This feature brings us two immediate benefits:(i) TJFast typically access much fewer elements than algorithms based on region encoding; and (ii) TJFast can efficiently process queries with wildcards in internal nodes. Our contributions in this paper can be summarized as follows:

- We propose an enhanced *Dewey ID* labeling scheme by incorporating element-name (i.e. node-type) information. Our approach is based on using *modulo* function and a *finite state transducer*(FST) to derive the element names along a path.

<sup>1</sup>Note that even if *b* and *c* are descendants of *a* and their level difference with *a* is 2, *b* and *c* may not be query answers, as they do not share the common parent.

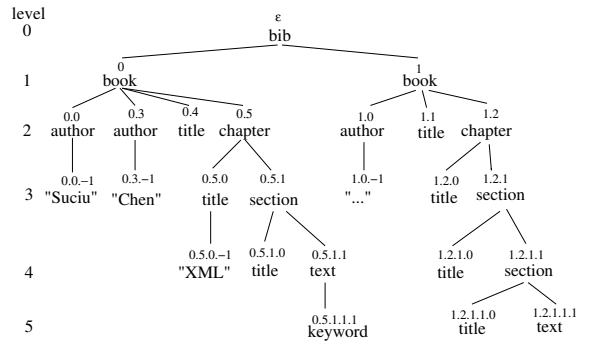


Figure 1: An XML tree with *extended Dewey* labels

- We develop a novel holistic twig join algorithm, called TJFast. When there are only A-D relationships between branching nodes and their children, TJFast is I/O optimal among all sequential algorithms that read the entire input. In other words, the optimality of TJFast allows the existence of P-C relationships between non-branching nodes and the children.
- We perform a comprehensive experiment to demonstrate the benefits of our algorithms over previous approaches.

**Organization** The rest of the paper proceeds as follows. We first discuss preliminaries in Section 2. The *extended Dewey* labeling scheme is presented in Section 3. We present TJFast algorithm in Section 4. Section 5 is dedicated to the related work. We present the experimental results in Section 6 and conclude this paper in Section 7.

## 2 Preliminaries

### 2.1 Data model and XML twig pattern

We model XML documents as *ordered* trees. 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 “\*”. The query twig pattern edges are either parent-child or ancestor-descendant edges. For convenience, 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.

Given a query twig pattern *Q* and an XML document *D*, a match of *Q* in *D* is identified by a mapping from the nodes in *Q* to the elements in *D*, such that: (i) the query node predicates are satisfied by the corresponding database elements, wherein wildcard “\*” can match any single tag; and (ii) the parent-child and ancestor-descendant relationships between query nodes are satisfied by the corresponding database elements. The answer to query *Q* with *n* nodes can be

represented as a list of  $n$ -ary tuples, where each tuple  $(q_1, \dots, q_n)$  consists of the database elements that identify a distinct match of  $Q$  in  $D$ .

## 2.2 Dewey ID labeling scheme

Tatarinov et al.[22] propose *Dewey ID* labeling scheme to present the position of an element occurrence in an XML document. In *Dewey ID*, each element is presented by a vector: (i) the root is labeled by a empty string  $\varepsilon$ ; (ii) for a non-root element  $u$ ,  $label(u) = label(s).x$ , where  $u$  is the  $x$ -th child of  $s$ . *Dewey ID* supports efficient evaluation of structural relationships between elements. That is, element  $u$  is an ancestor of element  $s$  if and only if  $label(u)$  is a prefix of  $label(s)$ .

*Dewey ID* has a nice property: one can derive the ancestors of an element from its label alone. For example, suppose element  $u$  is labeled “1.2.3.4”, then the parent of  $u$  is “1.2.3” and the grandparent is “1.2” and so on. With the knowledge of this property, we further consider that if the names of all ancestors of  $u$  can be derived from  $label(u)$  alone, then XML path pattern matching can be directly reduced to string matching. For example, if we know that the label “1.2.3.4” presents the path “a/b/c/d”, then it is quite straightforward to identify whether the element matches a path pattern (e.g. “//c/d”). Inspired by this observation, we develop an *extended Dewey ID* labeling scheme which provides an *extraordinary* chance for us to design a new algorithm to match XML path (and twig) pattern.

## 3 Extended Dewey and FST

In this section, we aim at extending *Dewey ID* labeling scheme to incorporate the element-name information. A straightforward way is to use some bits to present the element-name sequence with number presentation, followed by the *original Dewey* label. The advantage of this approach is simple and easy to implement. However, as shown in our experiments in Section 6, this method faces the problem of the large label size. In the following, we will propose a more concise scheme to solve this problem. In particular, we first *encode* the names of elements along a path into a single Dewey label. Then we present a *Finite State Transducer*(FST) to *decode* element names from this label. For simplicity, we focus the discussion on a single document. The labeling scheme can be easily extended to multiple documents by introducing document ID information.

### 3.1 Extended Dewey

The intuition of our method is to use *modulo* function to create a mapping from an integer to an element name, such that given a sequence of integers, we can convert it into the sequence of element names.

```

<!ELEMENT bib (book*)>
<!ELEMENT book ( author+, title, chapter* ) >
<!ELEMENT author (#PCDATA)>
<!ELEMENT title (#PCDATA)>
<!ELEMENT chapter (title, section*)>
<!ELEMENT section (title, (text | section)*)>
<!ELEMENT text (#PCDATA | bold | keyword | emph )*>
<!ELEMENT bold (#PCDATA | bold | keyword | emph )*>
<!ELEMENT keyword (#PCDATA | bold | keyword | emph )*>
<!ELEMENT emph (#PCDATA | bold | keyword | emph )*>

```

Figure 2: DTD for XML document in Fig 1

In the *extended Dewey*, we need to know a little additional schema information, which we call a *child names clue*. In particular, given any tag  $t$  in a document, the *child names clue* is all (distinct) names of children of  $t$ . This clue is easily derived from DTD, XML schema or other schema constraint. For example, consider the DTD in Figure 2; the tag of all children of *bib* is only *book* and the tags of all children of *book* are *author*, *title* and *chapter*. Note that even in the case when DTD and XML schema are unavailable, our method is still effective, but we need to scan the document once to get the necessary *child names clue* before labeling the XML document.

Let us use  $CT(t) = \{t_0, t_1, \dots, t_{n-1}\}$  to denote the *child names clue* of tag  $t$ . Suppose there is an ordering for tags in  $CT(t)$ , where the particular ordering is not important. For example, in Fig 3,  $CT(book) = \{author, title, chapter\}$ . Using child names clues, we may easily create a mapping from an integer to an element name. Suppose  $CT(t) = \{t_0, t_1, \dots, t_{n-1}\}$ , for any element  $e_i$  with name  $t_i$ , we assign an integer  $x_i$  to  $e_i$  such that  $x_i \bmod n = i$ . Thus, according to the value of  $x_i$ , it is easy to derive its element name. For example,  $CT(book) = \{author, title, chapter\}$ . Suppose  $e_i$  is a child element of *book* and  $x_i = 8$ , then we see that the name of  $e_i$  is *chapter*, because  $x_i \bmod 3 = 2$ . In the following, we extend this intuition and describe the construction of *extended Dewey* labels.

The *extended Dewey* label of each element can be efficiently generated by a *depth-first* traversal of the XML tree. Each *extended Dewey* label is presented as a vector of integers. We use  $label(u)$  to denote the *extended Dewey* label of element  $u$ . For each  $u$ ,  $label(u)$  is defined as  $label(s).x$ , where  $s$  is the parent of  $u$ . The computation method of integer  $x$  in *extended Dewey* is a little more involved than that in the *original Dewey*. In particular, for any element  $u$  with parent  $s$  in an XML tree,

- (1) if  $u$  is a text value, then  $x = -1$ ;
- (2) otherwise, assume that the element name of  $u$  is the  $k$ -th tag in  $CT(t_s)$  ( $k=0,1,\dots,n-1$ ), where  $t_s$  denotes the tag of element  $s$ .
  - (2.1) if  $u$  is the first child of  $s$ , then  $x = k$ ;
  - (2.2) otherwise assume that the last component of the label of the left sibling of  $u$  is  $y$  (at this point, the left sibling of  $u$  has been labeled), then



cestors of an element and therefore has not Properties 1 and 2. The original Dewey labeling scheme has Properties 2 to 5, but not Property 1. The first property is unique for *extended Dewey*. Note that Property 1 and 2 are of paramount importance, since they provide us an *extraordinary* chance to efficiently process XML path (and twig) queries. For example, given a path query  $a/b/c/d$ , according to the **Ancestor Name and Label Vision Properties**, we only need to read the labels of  $d$  to answer this query, which will significantly reduce I/O cost of previous algorithms based on region encoding. In the next section, we will use *extended Dewey* labels to design a novel and efficient holistic twig join algorithm, which utilizes the above five properties.

## 4 Twig Pattern Matching

### 4.1 Path matching algorithm

It is quite straightforward to evaluate a query path pattern in our approach. According to the Ancestor Name Vision and Label Properties, we *only need to scan the elements whose tags appear in leaf node of query*. For each visited element, we first use FST to reveal the element names along the whole path, and then perform string matching against it. As a result, we evaluate the path pattern efficiently by scanning the input list once and ensure that each output solution is our desired final answer.

When path queries contain only parent-child relationships within the path, the string-matching can be processed very efficiently by simply comparing element names. When path queries contain ancestor-descendant relationships or wildcards “\*”, the queries can be processed by string-matching with *don't care* symbols. There are a rich set of algorithms on efficient string processing with *don't care* symbols. (e.g. [18] and [9]).

It is worth noting that the I/O cost of our approach is typically much smaller than that of previous algorithms for path pattern matching (e.g. PathStack [3]), for we only scan labels for the query *leaf* node, while they need to scan elements for *all* query nodes.

### 4.2 Twig matching algorithm

This section presents a holistic twig pattern join algorithm, called TJFast. We will first introduce some data structures and notations.

#### 4.2.1 Data Structures and Notations

Let  $q$  denote a twig pattern and  $p_n$  denote a path pattern from the *root* to the node  $n \in q$ . In our algorithms, we make use of the following query node operations: `isleaf`: Node  $\rightarrow$  Bool; `isBranching`: Node  $\rightarrow$  Bool; `leafNodes`: Node  $\rightarrow$  {Node}; `directBranchingOrLeafNodes`: Node  $\rightarrow$  {Node}. `leafNodes( $n$ )` returns the set of leaf

nodes in the twig rooted with  $n$ . `directBranchingOrLeafNodes( $n$ )` (for short, `dbl( $n$ )`) returns the set of all branching nodes  $b$  and leaf nodes  $f$  in the twig rooted with  $n$  such that in the path from  $n$  to  $b$  or  $f$  (excluding  $n, b$  or  $f$ ) there is no branching nodes. For example, in the query Q1 of Fig 4, `dbl( $a$ )`={ $b, c$ } and `dbl( $c$ )`={ $f, g$ }.

Associated with each leaf node  $f$  in a query twig pattern there is a stream  $T_f$ . The stream contains *extended Dewey* labels of elements that match the node type  $f$ . The elements in the stream are sorted by the ascending lexicographical order. For example, “1.2” precedes “1.3” and “1.3” precedes “1.3.1”. The operations over a stream  $T_f$  include `current( $T_f$ )`, `advance( $T_f$ )` and `eof( $T_f$ )`. The function `current( $T_f$ )` returns the *extended Dewey* label of the current element in the stream  $T_f$ . The function `advance( $T_f$ )` updates the current element of the stream  $T_f$  to be its next element. The function `eof( $T_f$ )` tests whether we are in the end of the stream  $T_f$ . We make use of two self-explanatory operations over elements in the document: `ancestors( $e$ )` and `descendants( $e$ )`, which return the ancestors and descendants of  $e$ , respectively (both including  $e$ ).

Algorithm TJFast keeps a data structure during execution: a set  $S_b$  for each branching node  $b$ . Each two elements in set  $S_b$  have an *ancestor-descendant* or *parent-child* relationship. So the maximal size of  $S_b$  is *no more than the length of the longest path* in the document. Each element cached in sets likely participates in query answers. Set  $S_b$  is initially empty.

#### 4.2.2 TJFast

Algorithm TJFast, which computes answers to a query twig pattern  $q$ , is presented in Algorithm 1. TJFast operates in two phases. In the first phase (line 1-9), some solutions to individual root-leaf path patterns are computed. In the second phase (line 10), these solutions are merge-joined to compute the answers to the query twig pattern.

Given the *extended Dewey* label of an element, according to the **Ancestor Name Vision** property, it is easy to check whether its path matches the individual root-leaf path pattern. Thus, the key issue of TJFast is to determine whether a path solution can contribute to the solutions for the whole twig. In the optimal case, we only output the path solution that is merge-joinable to at least one solution of other root-leaf paths. Intuitively, if two path solutions can be merged, the necessary condition is that they have the common element to match the *branching* query node. For example, consider a simple query  $a[./b]/c$  and two path solution  $(a_1, b_1)$  and  $(a_2, c_1)$ . Observe that two solutions can be merged only if  $a_1 = a_2$ . Therefore, in TJFast, in order to determine whether a path solution contributes to final answers, we try to find the most likely elements that match branching nodes  $b$  and store them in the corresponding set  $S_b$ .

---

**Algorithm 1** TJFast

---

```
1: for each  $f \in \text{leafNodes}(\text{root})$ 
2:   locateMatchedLabel( $f$ )
3: endfor
4: while ( $\neg \text{end}(\text{root})$ ) do
5:    $f_{act} = \text{getNext}(\text{topBranchingNode})$ 
6:   outputSolutions( $f_{act}$ )
7:   advance( $T_{f_{act}}$ )
8:   locateMatchedLabel( $f_{act}$ )
9: end while
10: mergeAllPathSolutions()
```

Procedure locateMatchedLabel( $f$ )

/\* Assume that the path from the root to element  $\text{get}(T_f)$  is  $n_1/n_2/\dots/n_k$  and  $p_f$  denotes the path pattern from the root to leaf node  $f$  \*/

```
1: while  $\neg((n_1/n_2/\dots/n_k \text{ matches pattern } p_f) \wedge (n_k \text{ matches } f))$  do
2:   advance( $T_f$ )
3: end while
```

Function  $\text{end}(n)$

```
1: Return  $\forall f \in \text{leafNodes}(n) \rightarrow \text{eof}(T_f)$ 
```

Procedure outputSolutions( $f$ )

```
1: Output path solutions of  $\text{current}(T_f)$  to pattern  $p_f$  such that in each solution  $s$ ,  $\forall e \in s$ : (element  $e$  matches a branching node  $b \rightarrow e \in S_b$ )
```

---

It is not difficult to understand the main procedure of TJFast (see Algorithm 1). In line 1-3, for each stream, we use Procedure locateMatchedLabel to locate the first element whose path matches the individual root-leaf path pattern. In line 5, we identify the next stream  $T_{f_{act}}$  to be processed by using  $\text{getNext}(\text{topBranchingNode})$  algorithm, where  $\text{topBranchingNode}$  is defined as the branching node that is an ancestor of all other branching nodes (if any). In line 6, we output some path matching solutions in which each element that match any branching node  $b$  can be found in the corresponding set  $S_b$ . We advance  $T_{f_{act}}$  in line 7 and locate the next matching element in line 8.<sup>2</sup>

Algorithm  $\text{getNext}$  (see Algorithm 2) is the core function called in TJFast, in which we accomplish two tasks. The first is to identify the next stream to process; and the second is to update the sets  $S_b$  associated with branching nodes  $b$ , discussed as follows.

For the first task to identify the next processed stream, Algorithm  $\text{getNext}(n)$  returns a query leaf node  $f$  according to the following recursive criteria (i) if  $n$  is a leaf node, return  $n$  (line 2); else (ii)  $n$  is a branching node, then for each node  $n_i \in \text{dbl}(n)$ , (1)

---

<sup>2</sup>Note that the second condition “ $n_k$  matches  $f$ ” in line 1 of locateMatchedLabel is necessary, which avoids outputting duplicate solutions. For example, consider the element  $e$  (with tag name  $b$ ) with the path “ $a_1/b_1/c_1/b_2$ ” and the path query “ $a/b$ ”. “ $a_1/b_1/c_1/b_2$ ” can match “ $a/b$ ”, but this solution has been output by another element ends with  $b_1$ .

---

**Algorithm 2** getNext( $n$ )

---

```
1: if ( $\text{isLeaf}(n)$ ) then
2:   return  $n$ 
3: else
4:   for each  $n_i \in \text{dbl}(n)$  do
5:      $f_i = \text{getNext}(n_i)$ 
6:     if ( $\text{isBranching}(n_i) \wedge \text{empty}(S_{n_i})$ )
7:       return  $f_i$ 
8:      $e_i = \max\{p \mid p \in MB(n_i, n)\}$ 
9:   end for
10:   $\max = \max_{arg_i}\{e_i\}$ 
11:   $\min = \min_{arg_i}\{e_i\}$ 
12:  for each  $n_i \in \text{dbl}(n)$  do
13:    if ( $\forall e \in MB(n_i, n) : e \notin \text{ancestors}(e_{\max})$ )
14:      return  $f_i$ ;
15:    endif
16:  end for
17:  for each  $e \in MB(n_{\min}, n)$ 
18:    if ( $e \in \text{ancestors}(e_{\max})$ ) updateSet( $S_n, e$ )
19:  end for
20:  return  $f_{\min}$ 
21: end if
```

Function MB( $n, b$ )

```
1: if ( $\text{isBranching}(n)$ ) then
2:   Let  $e$  be the maximal element in set  $S_n$ 
3: else
4:   Let  $e = \text{current}(T_n)$ 
5: end if
6: Return a set of element  $a$  that is an ancestor of  $e$  such that  $a$  can match node  $b$  in the path solution of  $e$  to path pattern  $p_n$ 
```

Procedure clearSet( $S, e$ )

```
1: Delete any element  $a$  in the set  $S$  such that  $a \notin \text{ancestors}(e)$  and  $a \notin \text{descendants}(e)$ 
```

Procedure updateSet( $S, e$ )

```
1: clearSet( $S, e$ )
2: Add  $e$  to set  $S$ 
```

---

if the current elements cannot form a match for the subtree rooted with  $n_i$ , we immediately return  $f_i$  (line 7); (2) if the current element from stream  $T_{f_i}$  does not participate in the solution involving in the future elements in other streams, we return  $f_i$  (line 14); (3) otherwise we return  $f_{\min}$  such that the current element  $e_{\min}$  has the minimal label in all  $e_i$  by lexicographical order (line 20).

For the second task, we update set  $e_b$ . This operation is important, since the elements in  $e_b$  decides which path solution can be output in Procedure  $\text{outputSolutions}$ . In line 18 of Algorithm 2, before an element  $e_b$  is inserted to the set  $S_b$ , we ensure that  $e_b$  is an ancestor of (or equals) each other element  $e_{b_i}$  to match node  $b$  in the corresponding path solutions.

EXAMPLE 4.1 Consider Q1 and Doc1 in Fig 4(a-b). A subscript is added to each element in the order of

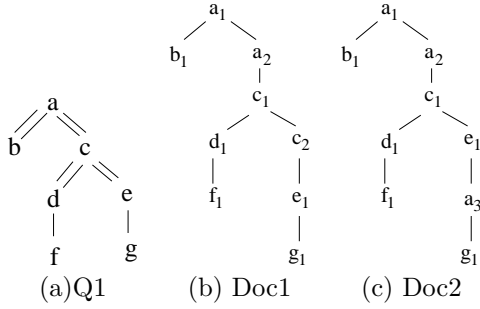


Figure 4: Example twig query and documents

pre-order traversal for easy reference. There are three input streams  $T_b$ ,  $T_f$  and  $T_g$ . Initially,  $getNext(a)$  recursively calls  $getNext(b)$  and  $getNext(c)$  (for  $b, c \in dbl(a)$  in  $Q1$ ). Since  $b$  is a leaf node in  $Q1$ ,  $getNext(b)=b$ . Observe that  $MB(f,c)=\{c_1\}$  and  $MB(g,c)=\{c_1, c_2\}$ , So  $e_{max} = g$  and  $e_{min} = f$  in line 10 and 11 of Algorithm 2. In line 18,  $c_1$  is inserted to set  $S_c$ . Then,  $getNext(c)=f$ . Subsequently,  $a_1$  is inserted to  $S_a$  and  $getNext(a)=b$ . Finally path solutions  $(a_1, b_1), (a_1, c_1, d_1, f_1)$  and  $(a_1, c_1, e_1, g_1)$  are output and merged. Note that although  $(a_1, c_2, e_1, g_1)$  matches the individual path pattern  $a/c/e/g$ , it is not output for  $c_2 \notin S_c$ .  $\square$

Note that the second phase(line 10 of Algorithm 1) of TJFast can be performed efficiently, only when the intermediate path solutions are output in sorted order. To achieve this purpose, we would need to “block” some answers. The details of how to achieve this naturally in the scenario of TJFast can be found in [15] and are omitted here for reason of space.

### 4.3 Analysis of TJFast

Next, we first show the correctness of TJFast and then analyze its complexity.

**Lemma 1.** *In Procedure clearSet of Algorithm TJFast, any element  $e$  that is deleted from set  $S_b$  does not participate in any new solution.*

**Lemma 2.** *In line 18 of Function getNext, if element  $e \notin ancestors(e_{max})$  and  $e \notin S_n$ , then  $e$  is guaranteed to not involve in any final solution.*

Lemma 1 shows that any element deleted from sets does not participate in new solutions, so the deletion is safe. Lemma 2 shows that for any element  $e$  that matches a branching node, if  $e$  participates in any final answer, then  $e$  occurs in the corresponding set. Thus the insertion is complete. The two lemmas are important to establish the correctness of the following theorem.

**Theorem 1.** *Given a twig query  $Q$  and an XML database  $D$ , Algorithm TJFast correctly returns all the answers for  $Q$  on  $D$ .*

While the correctness holds for any given query, the I/O optimality holds only for the case where there are only ancestor-descendant relationships between branching nodes and their children.

**Theorem 2.** *Consider an XML database  $D$  and a twig query  $Q$  with only ancestor-descendant relationships between branching nodes and their children. The worst case I/O complexity of TJFast is linear in the sum of the sizes of input and output lists. The worst-case space complexity is  $O(d^2 * |b| + d * |f|)$ , where  $|f|$  is the number of leaf nodes in  $q$ ,  $|b|$  is the number of branching nodes in  $q$  and  $d$  is the length of the longest label in the input lists.*

**PROOF:[sketch]** We first prove the I/O optimality. The following observation is important to prove the optimality of TJFast: if all branching edges are only ancestor-descendant relationships, then in line 18 of getNext, since  $e \in ancestors(e_{max})$ ,  $e \in MB(n_i, n)$  for each  $n_i \in dbl(n)$ . That is,  $e$  is guaranteed to be a common element in each current path solution. Note that we only output path solutions, in which elements that match branching nodes occur in the corresponding set(line 6 of Algorithm 1). Therefore, each intermediate path solution output in TJFast is guaranteed to contribute to final results when the query contains only ancestor-descendant relationships in branching edges.

As for space complexity, our result is based on the observation that in the worst case, the number of elements in branching node set  $S_b$  is at most  $d$ , where  $d$  is the length of the longest label in the input lists. Considering each extended Dewey label repeats its prefix, the total space complexity of  $S_b$  is  $O(d^2)$ .  $\square$

Theorem 2 holds only for query with ancestor-descendant relationships to connect branching nodes. Unfortunately, in the case where the query contains parent-child relationships between branching nodes and their children, Algorithm TJFast is no longer guaranteed to be I/O optimal. For example, consider a query  $a[./b]/c$  and a data tree consisting of  $a_1$ , with children(in order)  $b_1, a_2, c_2$ , such that  $a_2$  has children  $b_2, c_1$ . There are two streams  $T_b, T_c$  in TJFast and their first elements are  $b_1$  and  $c_1$  respectively. In this case,  $b_1$  and  $c_1$  are “locked” simultaneously, because we cannot advance any stream before knowing if it participates in a solution. Thus, optimality can no longer be guaranteed.

### 4.4 Comparison among TJFast, TwigStack and TwigStackList

In this section, we use the following example to illustrate the advantages of TJFast over TwigStack and TwigStackList.

**EXAMPLE 4.2** Consider the query and data tree Doc2 in Fig 4(a) and (c). There are three input streams  $T_b, T_f$  and  $T_g$  in TJFast. Initially, the current elements are  $b_1, f_1$  and  $g_1$ . TJFast does not insert

$c_1$  to set  $S_c$ , since by reading the label of  $g_1$  alone, we immediately identify that  $g_1$  does not contribute to query answers (for  $a_1/a_2/c_1/e_1/a_3/g_1$  does not match  $a//c//e/g$ ). In contrast, TwigStack pushes  $c_1$  to stack  $S_c$  and outputs two “useless” intermediate path solution  $\langle a_1, b_1 \rangle$  and  $\langle a_1, c_1, d_1, f_1 \rangle$ . The behavior of TwigStack is also reasonable because based on *region coding* of  $g_1$ , one cannot decide whether  $g_1$  has the parent tagged with  $e$ . But based on *extended Dewey*, one can easily identify that the parent of  $g_1$  is tagged with  $a$  rather than  $e$ . This example shows the benefit of *extended Dewey* labeling scheme on efficient processing of XML twig pattern matching.

Compared to TwigStack, TwigStackList looks more “clever”. In the above example, TwigStackList does not hastily push  $c_1$  to stack, but first checks the parent-child relationship between  $e_1$  and  $g_1$ . Then they find that  $e_1$  is not the parent of  $g_1$ . Then TwigStackList caches  $e_1$  in a list and reads more elements in  $T_e$ . In this simple case,  $e_1$  is the only element in stream  $T_e$ . So unlike TwigStack, TwigStackList does not output any useless intermediate results. Compared to TJFast, TwigStackList is also I/O optimal in this example, but TwigStackList needs to read more elements from all non-leaf node streams and its processing will be very complicated when  $g_1$  has more than one ancestor tagged with  $e$ . (More examples about TwigStackList can be found in [13]) □

## 5 Related work

**Labeling schemes** *Dewey ID* labeling scheme comes from the work of Tatarinov et al.[22] to represent XML order in the relational data model, and to show how this labeling scheme can be used to preserve document order during XML query processing. O’Neil et al.[17] introduced a variation of prefix labeling scheme called ORDPATH. Unlike our *extended Dewey*, the main goal of ORDPATH is to gracefully handle insertion of XML nodes in the database.

The *region encoding* is considered as the work of Consens and Milo[8], who discuss a fragment of PAT text searching operators for indexing text database. Then Zhang et al.[27] introduced it to XML query processing using inverted list. Recently, many researchers ([4],[21],[25]) have begun to design dynamic XML labeling schemes to handle data updates.

**Twig join algorithms** Al-Khalifa et al.[1] started the stack-based algorithms for XML structural joins. N. Bruno et al. [3] proposed a holistic twig join algorithm, namely TwigStack. Lu et al.[13] proposed TwigStackList, which identifies a larger optimal query class than TwigStack. Lu et al.[14] also researched how to answer an *ordered* twig pattern based on region encoding. Chen et al.[6] proposed an algorithm iTwigJoin, which is still based on region encoding, but work with different data partition strategies (e.g. Tag+Level and Prefix Path Streaming).

Jiang et al. [11] proposed a general algorithm called TSGeneric+ based on indexes built on element labels. Their method can skip elements and achieve sub-linear performance for selective queries. But for evaluating queries with parent-child relationships, TSGeneric+ may still output many “useless” intermediate results like TwigStack. Jiang et al.[10] also studied the problem of processing queries with OR predicates. BLAS by Chen et al. [7] proposed a bi-labelling scheme: D-Label and P-Label for accelerating *parent-child* relationship processing. Their method decomposes a twig pattern into several *parent-child* path queries and then merges the results.

Yang et al. [26] proposed the idea of combining path index table and Dewey labels.<sup>3</sup> Similar to our TJFast, to answer a twig query, their method also can reduce I/O cost by accessing only the labels of leaf query nodes. But unlike TJFast, their algorithm did not fully exploit the nice properties of *Dewey* labels and only modified one procedure in TSGeneric+. So similar to TSGeneric+, their algorithm is still not efficient for processing queries with parent-child relationships.

ViST and PRiX ([24],[19]) transform both XML data and queries into sequences and answer XML queries through subsequence matching. While their methods avoid join operations in query processing, to eliminate false alarm and false dismissal, they resort to post-processing(for false alarm) and multiple isomorphism queries processing(for false dismissal[23]), both of which are time consuming.

## 6 Experimental study

### 6.1 Experimental setup

We implemented four XML twig join algorithms: TJFast, TwigStack, TwigStackList and iTwigJoin in JDK 1.4 using the file system for storage. Only TJFast is based on *extended Dewey* labeling scheme, and the other three use *region encoding*.

The reason that we chose these three algorithms is that they are efficient for different applications. TwigStack[3] is very efficient when query contains only ancestor-descendant relationships. TwigStackList[13] is efficient on answering queries with parent-child relationships. Finally, unlike the above two algorithms, which partition elements based on their tags, iTwigJoin[6] is a general twig join algorithm, which can be used on different data partitioning approaches. [6] researched two new data partitions: *tag+level* and *prefix path streaming* (PPS). Such *refined* data partitioning strategies enable iTwigJoin to reduce I/O cost by pruning irrelevant data streams.

All experiments were run on a 1.7G Pentium IV processor running Windows XP with 768MB of main

<sup>3</sup>Note that our work was developed independently of and differs considerably from [26].

Table 1: XML Data Sets (XM: XMark,TB:TreeBank)

	XM	Random	DBLP	TB
Data size(MB)	582	90	130	82
Nodes(million)	8	5.1	3.3	2.4
Max/Avg depth	12/5	10/5.1	6/2.9	36/7.8

Table 2: Labels size (XM: XMark,TB:TreeBank)

	XM	Random	DBLP	TB
Original Dewey(MB)	56.2	36.1	18.1	22.8
Region coding(MB)	71.9	45.2	21.6	23.3
Naive extension(MB)	92.9	55.8	27.7	41.9
Extended Dewey(MB)	72.6	43.3	19.5	28.7

memory and 2GB of disk space. We used four different datasets, including two synthetic and two real datasets. The first synthetic dataset is the well-known XMark benchmark data (with factor 5). The second is a random data set with ten distinct labels (namely  $A_1, A_2, \dots, A_{10}$ ). The node labels in the tree were uniformly distributed. The two real datasets are DBLP and TreeBank[16]<sup>4</sup>. We chose these two datasets since they have different characteristics: DBLP is a shallow and wide document, but TreeBank has very deep recursive structure. Table 1 summarizes the characteristics of the four datasets.

In our experiments, the *extended Dewey* labels are not stored by the dotted-decimal strings displayed (e.g. “1.2.3.4”), but rather a compressed binary representation. In particular, we used UTF-8 encoding as an efficient way to present the integer value, which was proposed by Tatarinov et al. [22]. Our experimental results show that compared to the naive implementation, where each integer value is presented as a fixed number of bytes, the UTF-8 encoding can save about 50% space cost.

## 6.2 Experimental results

### 6.2.1 Labels size

We compared the labels sizes of four labeling schemes in Table 2. Our first conclusion is that the size of the *naive extension*, which directly presents the element-name sequence in number presentation ahead of the *original Dewey* labels, is generally larger than that of our *extended Dewey* labeling scheme. Our second conclusion is that when the document tree is shallow and wide (i.e. DBLP), the size of *extended Dewey* is smaller than that of *region encoding*. But when the document tree is deep (i.e. TreeBank), the size of *region encoding* is smaller. This is because *extended Dewey* is a variation of prefix labeling scheme, whose

<sup>4</sup>Since there is no DTD available for TreeBank and random data, we first scan this document once to get the *child names clue* of each tag.

size is closely related to the average depth of documents. Our third conclusion is that the size of *extended Dewey* is about 10%-30% more than that of *original Dewey*. As we will show in our experiments, it is worth using this additional space-overhead, since it significantly improves the performance of XML twig pattern matching.

### 6.2.2 Path Queries

We next compare our algorithm TJFast with the previous PathStack[3] to match path queries without branching nodes. For this purpose we used XMark benchmark data and four path queries<sup>5</sup> shown in Table 3. Figure 5 compares two algorithms in terms of the number of elements read, the size of disk files scanned and execution time.

An immediate observation from the figures is that TJFast is more efficient than PathStack. In particular, PathStack could perform 400% more disk I/Os than those required by TJFast (e.g.  $PQ_2$ ).

In order to research the effect of query path length on TJFast and PathStack, we then used the random data set consisting of ten distinct labels  $A_1, A_2, \dots, A_{10}$ , and issue path queries of different lengths such as  $A_1/A_2/\dots/A_{10}$ . Figure 6 shows the execution times of both techniques, as well as the number of elements read and the size of disk files. Clearly, TJFast results in considerably better performance than PathStack. The performance of PathStack degrades significantly with the increase of the path length, but that of TJFast is almost not affected at all, as TJFast only scan data associated with the leaf node.

Table 3: Path Queries on XMark data

	Path Queries
$PQ_1$	/site/closed_auctions/closed_auction/price
$PQ_2$	/site/regions//item /location
$PQ_3$	/site/people/person/gender
$PQ_4$	/site/open_auctions/open_auction/reserve

### 6.2.3 Twig Queries

Table 4: Twig Queries on DBLP and TreeBank(TB)

	Data	Type	Twig Queries
$TQ_1$	DBLP	1	//inproceedings//title[./i]//sup
$TQ_2$	DBLP	1	//article[./sup]//title//sub
$TQ_3$	TB	2	/S[./VP/IN]//NP
$TQ_4$	TB	3	/S/VP/PP[IN]/NP/VBN
$TQ_5$	TB	4	//VP[DT]//PRP_DOLLAR_

We now focus on twig queries, and compare four holistic twig join algorithms TwigStack, TwigStack-

<sup>5</sup>We chose these queries according to XMark benchmark queries in [20].

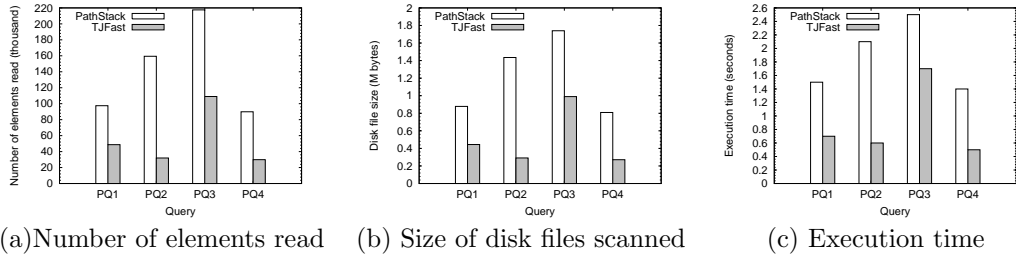


Figure 5: PathStack versus TjFast using XMark data

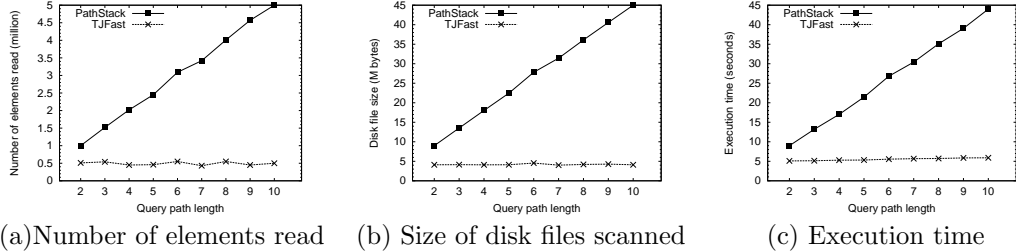


Figure 6: PathStack versus TjFast using random data

List, iTwigJoin and TjFast. We tested several XML queries on DBLP and TreeBank data (see Table 4)<sup>6</sup>. These queries have different twig structures and combinations of parent-child and ancestor-descendant relationships. In particular, queries TQ1 and TQ2 contain only ancestor-descendant relationships, while TQ4 contains only parent-child relationships. TQ3 contains only ancestor-descendant relationships between the branching node and its children, while TQ5 contains a branching node with both parent-child and ancestor-descendant relationships.

**TjFast vs. TwigStack** We first compare the performance between TjFast and TwigStack. From Figure 7 and 8, we see that TjFast outperforms TwigStack for all queries. We now analyze the query performance under two scenarios namely *the cost of disk access* and *the size of intermediate results*.

*Cost of disk access* Figures 7(a) and 8(a) show that TjFast read far fewer elements than TwigStack. For example, for TQ1, TwigStack read 442167 elements, but TjFast read only 2380 elements (over two orders of magnitude). This huge gap results from the fact that TwigStack scans the elements for *all* the queries nodes, but TjFast scans only elements for *leaf* nodes.

*Size of intermediate results* Table 5 shows the number of intermediate path solutions output by different algorithms. The last column is the number of intermediate solutions that contribute to the final answers. An immediate observation is that TwigStack outputs many “*useless*” path solutions to queries with parent-child edges. For example, for TQ<sub>3</sub>, TwigStack produced 702391 intermediate paths, of which only 22565

are useful. More than 95% intermediate solutions output by TwigStack are “*useless*” to the final answers. In contrast, TjFast is optimal for query TQ<sub>3</sub> since the number of paths produced by TjFast is equal to the number of useful solutions.

Table 5: Number of intermediate path solutions

Query	TwigStack	TwigStackList	TjFast	Useful
TQ <sub>3</sub>	702391	22565	22565	22565
TQ <sub>4</sub>	2237	388	388	302
TQ <sub>5</sub>	10663	9	9	5

**TjFast vs. TwigStackList** From Fig. 7 and 8, TjFast also outperforms TwigStackList for all queries. This can be explained by the fact that TjFast reduces the I/O cost of TwigStackList by reading labels of only the *leaf* nodes.

When queries contain parent-child relationships between the branching node and its children (i.e. queries TQ<sub>4</sub>, TQ<sub>5</sub>), both TwigStackList and TjFast are sub-optimal. Their sub-optimality is evident from the observation that the number of intermediate path solutions by TwigStackList and TjFast is slightly larger than the number of useful solutions.

**TjFast vs. iTwigJoin** We now compare the performance between TjFast and iTwigJoin. iTwigJoin is based on region encoding, but it can be applied with different data partitioning strategies. Since [6] proposed two new data partitioning strategies (i.e. Tag+Level and PPS), we compare both variants with TjFast (labeled as iTwigJoin-TL and iTwigJoin-PPS, respectively).

Figure 9 and 10 compare the performance of

<sup>6</sup>We tried twig queries on XMark data. Those results are omitted due to space limitation and can be found in [15].

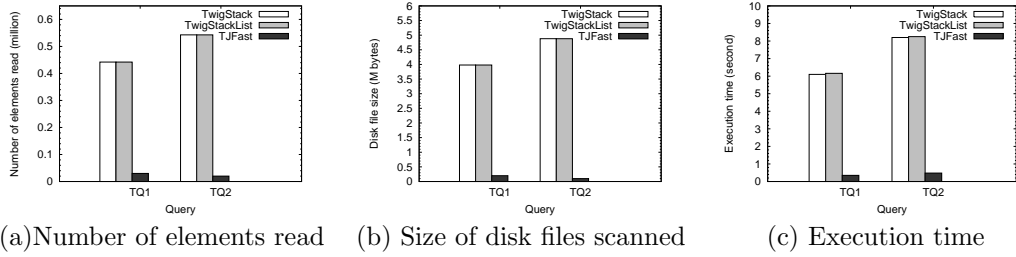


Figure 7: TwigStack, TwigStackList versus TJFast on DBLP

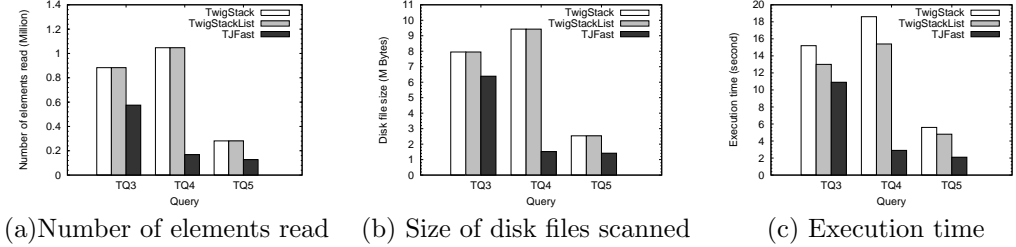


Figure 8: TwigStack, TwigStackList, TJFast on TreeBank

iTwigJoin-TL, iTwigJoin-PPS and TJFast on DBLP and TreeBank datasets. Since [6] has shown that PPS is not applicable to *deep* recursive data, for TreeBank, we only compared iTwigJoin-TL with TJFast. As shown from these results, we can see TJFast is again more efficient than iTwigJoin-TL and iTwigJoin-PPS for all queries. Although iTwigJoin uses the refined data partitioning strategies and scan less elements than TwigStack and TwigStackList, the number of elements processed by iTwigJoin is still more than that by TJFast.

## 6.2.4 Wildcard Queries

Finally, we tested two wildcard queries Q1://NP[./CD]\*/V and Q2://VP/\*[PP-8]/PP-7 on TreeBank dataset. Q1 is a twig query consisting of a wildcard in a non-branching node, but Q2 is a branching wildcard twig query. For Q1, all four algorithms can be applied. But the performance of TJFast is much better than the best algorithm based on region encoding<sup>7</sup>(0.9s vs. 7.2s). For Q2, the algorithms using region encoding are significantly *affected* by wildcards in branching nodes, as they do not know which elements can be used to match this wildcard. Since there is no DTD available for TreeBank data, a brute-force solution is to access all elements to answer this query. Clearly, this method is unacceptably slow. In contrast, the existence of wildcard in branching nodes *does not affect* TJFast, which takes only 0.3s to answer Q2. This shows that TJFast supports efficient processing of both non-branching as well as branching wildcard queries.

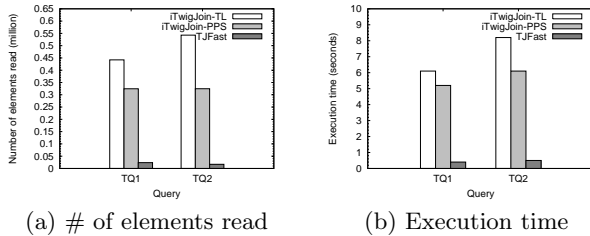


Figure 9: iTwigJoin, TJFast on DBLP

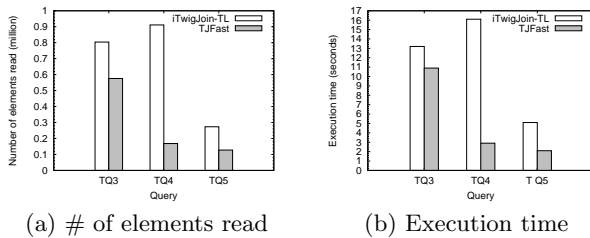


Figure 10: iTwigJoin, TJFast on TreeBank

**Summary** TJFast significantly outperforms TwigStack, TwigStackList and iTwigJoin under all settings (including shallow and deep documents, path and twig queries, branching and non-branching wildcards queries). The improvement is due to the facts that TJFast only scans labels for query *leaf* nodes. Algorithms based on region encoding are comparable to TJFast only when the number of elements associated with all *internal* query nodes is very small.

<sup>7</sup>In this case the best algorithm on region encoding is iTwigJoin-TL.

## 7 Conclusions and Future Work

XML twig pattern matching is a key issue for XML query processing. In this paper, we have proposed TJFast as an efficient algorithm to address this problem using a novel labeling scheme called *extended Dewey*. Although the idea of *original Dewey* is not new, extending it to efficiently process XML twig pattern matching is nontrivial. This is because based on the *original Dewey*, we cannot know the element names along a path. To answer a twig query, we need to access the labels of *all* query nodes. Considering the fact that prefix comparison is less efficient than integer comparison, the performance of algorithm with the *original Dewey* is usually worse than that with *region encoding*. However, owing to our extension, *extended Dewey* has the important property: **Ancestor Name Vision**. So TJFast only needs to access labels of *leaf* nodes to answer queries and significantly reduce I/O cost. Further, TJFast can efficiently evaluate queries with wildcards steps, which cannot be handled by algorithms with region encoding. As part of future work, we would like to improve *extended Dewey* to become an insert-friendly labeling scheme in the context of dynamic XML trees.

## 8 Acknowledgment

We would like to thank the anonymous reviewers of VLDB for their constructive and valuable comments. Furthermore, we thank Beverly Yang for bringing our attention to a related paper [26].

## References

- [1] S. Al-Khalifa, H. V. Jagadish, J. M. Patel, Y. Wu, N. Koudas, and D. Srivastava. Structural joins: A primitive for efficient XML query pattern matching. In *Proc. of ICDE Conference*, pages 141–152, 2002.
- [2] J.-M. Bremer and M. Gertz. An efficient XML node identification and indexing scheme. Technical Report CSE-2003-04, University of California at Davis, 2003.
- [3] N. Bruno, D. Srivastava, and N. Koudas. Holistic twig joins: optimal XML pattern matching. In *SIGMOD Conference*, pages 310–321, 2002.
- [4] B. Catania, B. C. Ooi, W. Wang, and X. Wang. Lazy xml updates: Laziness as a virtue of update and structural join efficiency. In *SIGMOD*, To appear 2005.
- [5] C. Y. Chan, W. Fan, and Y. Zeng. Taming XPath queries by minimizing wildcard steps. In *Proceeding of VLDB*, pages 156–167, 2004.
- [6] T. Chen, J. Lu, and T. Ling. On boosting holism in XML twig pattern matching using structural indexing techniques. In *SIGMOD To appear*, 2005.
- [7] Y. Chen, S. B. Davidson, and Y. Zheng. BLAS: An efficient XPath processing system. In *Proc. of SIGMOD*, pages 47–58, 2004.
- [8] M. P. Consens and T. Milo. Optimizing queries on files. In *SIGMOD*, pages 301–312, 1994.
- [9] G. H. Gonnet. The PAT text searching system. Technical report, University of Waterloo, 1987.
- [10] H. Jiang, H. Lu, and W. Wang. Efficient processing of XML twig queries with OR-predicates. In *Proc. of SIGMOD Conference*, pages 274–285, 2004.
- [11] H. Jiang, W. Wang, and H. Lu. Holistic twig joins on indexed XML documents. In *Proc. of VLDB*, pages 273–284, 2003.
- [12] Q. Li and B. Moon. Indexing and querying XML data for regular path expressions. In *Proc. of VLDB*, pages 361–370, 2001.
- [13] J. Lu, T. Chen, and T. W. Ling. Efficient processing of XML twig patterns with parent child edges: a look-ahead approach. In *CIKM*, pages 533–542, 2004.
- [14] J. Lu, T. Ling, T. Yu, C. Li, and W. Ni. Efficient processing of ordered XML twig pattern matching. In *DEXA To appear*, 2005.
- [15] J. Lu, T. W. Ling, C. Y. Chan, and T. Chen. From region encoding to extended dewey: On efficient processing of xml twig pattern matching. Technical report, TRA6/05 National university of Singapore, 2005.
- [16] U. of Washington XML Repository. <http://www.cs.washington.edu/research/xmldatasets/>.
- [17] P. O’Neil, E. O’Neil, S. Pal, I. Cseri, G. Schaller, and N. Westbury. ORDPATHs: Insert-friendly XML node labels. In *SIGMOD*, pages 903–908, 2004.
- [18] R. Y. Pinter. Efficient string matching with don’t care patterns. In *Combinatorial Algorithms on Words, NATO ASI Series*, volume 12, pages 11–29, 1985.
- [19] P. Rao and B. Moon. PRIX: Indexing and querying XML using prufer sequences. In *ICDE*, pages 288–300, 2004.
- [20] A. R. Schmidt et al. XMark an XML benchmark project. <http://monetdb.cwi.nl/xml/index.html>.
- [21] A. Silberstein, H. He, K. Yi, and J. Yang. Boxes: Efficient maintenance of order-based labeling for dynamic XML data. In *Proc. of ICDE.*, pages 285–296, 2005.
- [22] I. Tatarinov, S. Vlas, K. S. Beyer, J. Shanmugasundaram, E. J. Shekita, and C. Zhang. Storing and querying ordered XML using a relational database system. In *Proc. of SIGMOD*, pages 204–215, 2002.
- [23] H. Wang and X. Meng. On the sequencing of tree structures for XML indexing. In *ICDE*, pages 372–383, 2005.
- [24] H. Wang, S. Park, W. Fan, and P. S. Yu. ViST: A dynamic index method for querying XML data by tree structures. In *SIGMOD*, pages 110–121, 2003.
- [25] X. Wu, M. Lee, and W. Hsu. A prime number labeling scheme for dynamic ordered XML trees. In *Proc. of ICDE*, pages 66–78, 2004.
- [26] B. Yang, M. Fontoura, E. J. Shekita, S. Rajagopalan, and K. S. Beyer. Virtual cursors for XML joins. In *CIKM*, pages 523–532, 2004.
- [27] C. Zhang, J. F. Naughton, D. J. DeWitt, Q. Luo, and G. M. Lohman. On supporting containment queries in relational database management systems. In *Proc. of SIGMOD Conference*, pages 425–436, 2001.