

## An Efficient Algorithm for Tree Mapping in XML Databases

Yangjun Chen

University of Winnipeg, Winnipeg, Manitoba, Canada R3B 2E9

---

**Abstract:** In this article, we discuss an efficient algorithm for tree mapping problem in XML databases. Given a target tree  $T$  and a pattern tree  $Q$ , the algorithm can find all the embeddings of  $Q$  in  $T$  in  $O(|T||Q|)$  time while the existing approaches need exponential time in the worst case.

**Key words:** Tree mapping, XML databases, query evaluation, tree encoding

---

### INTRODUCTION

XML uses a tree-structured model for representing data. Queries in XML languages (such as XPath<sup>[1]</sup>, Xquery<sup>[2,3]</sup>, XML-QL<sup>[4]</sup> and Quilt<sup>[5,6]</sup>) also typically specify selection patterns as a kind of tree-structured relations. For instance, the XPath expression:

`book[title = 'Art of Programming']/author[fn = 'Donald' and ln = 'Knuth']`

matches *author* elements that (i) have a child subelement *fn* with content 'Donald', (ii) have a child subelement *ln* with content 'Knuth' and are descendants of *book* elements that have a child *title* subelement with content 'Art of Programming'. This expression can be represented as a tree structure as shown in Fig. 1.

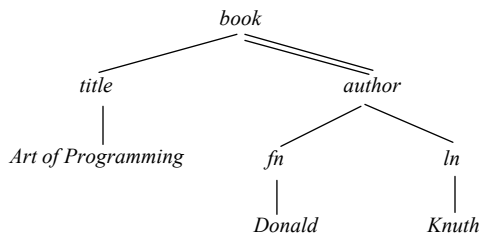


Fig. 1: A query tree

In this tree structure, a node  $v$  is labeled with an element name or a string value, denoted  $label(v)$ . In addition, there are two kinds of edges: child edges ( $c$ -edges) for parent-child relationships and descendant edges ( $d$ -edges) for ancestor-descendant relationships. A  $c$ -edge from node  $v$  to node  $u$  is denoted by  $v \rightarrow u$  in the text and represented by a single arc;  $u$  is called a  $c$ -child of  $v$ . A  $d$ -edge is denoted  $v \Rightarrow u$  in the text and represented by a double arc;  $u$  is called a  $d$ -child of  $v$ . Such a query is often called a twig pattern.

In any DAG (*directed acyclic graph*), a node  $u$  is said to be a descendant of a node  $v$  if there exists a path (sequence of edges) from  $v$  to  $u$ . In the case of a twig pattern, this path could consist of any sequence of  $c$ -edges and/or  $d$ -edges. Based on these concepts, the tree embedding can be defined as follows.

**Definition 1:** An embedding of a twig pattern  $Q$  into an XML document  $T$  is a mapping  $f: Q \rightarrow T$ , from the nodes of  $Q$  to the nodes of  $T$ , which satisfies the following conditions:

- i. Preserve node label: For each  $u \in Q$ ,  $u$  and  $f(u)$  are of the same label (or more generally,  $u$ 's predicate is satisfied by  $f(u)$ .)
- ii. Preserve  $c/d$ -child relationships: If  $u \rightarrow v$  in  $Q$ , then  $f(v)$  is a child of  $f(u)$  in  $T$ ; if  $u \Rightarrow v$  in  $Q$ , then  $f(v)$  is a descendant of  $f(u)$  in  $T$ .

If there exist a mapping from  $Q$  into  $T$ , we say,  $Q$  can be imbedded into  $T$ , or say,  $T$  contains  $Q$ .

Notice that an embedding could map several nodes of the query (of the same type) to the same node of the database. It also allows a tree mapped to a path. This definition is quite different from the tree matching defined in<sup>[7]</sup>.

There is much research on how to find such a mapping efficiently and all the proposed methods can be categorized into two groups. By the first group<sup>[2,8-17]</sup>, a tree pattern is typically decomposed into a set of binary relationships between pairs of nodes, such as parent-child and ancestor-descendant relations. Then, an index structure is used to find all the matching pairs that are joined together to form the final result. By the second group<sup>[18-23]</sup>, a query pattern is decomposed into a set of paths. The final result is constructed by joining all the matching paths together. For all these methods, the join operations involved require exponential time in the worst case. For example, if we decompose a twig pattern into paths to find all the matching paths from a

database, we need  $O(p^\lambda)$  time to join them together, where  $p$  is the largest length of a matching path and  $\lambda$  is the number of all such paths.

In this study, we proposed a new algorithm with no join operations involved. The algorithm runs in  $O(|T| \cdot Q_{leaf})$  time and  $O(T_{leaf} \cdot Q_{leaf})$  space, where  $T_{leaf}$  and  $Q_{leaf}$  represent the numbers of the leaf nodes in  $T$  and in  $Q$ , respectively.

### TREE ENCODING

To facilitate the checking of reachability (whether a node can be reached from another node through a path), a tree encoding is used<sup>[24]</sup>.

Consider a tree  $T$ . By traversing  $T$  in *preorder*, each node  $v$  will obtain a number  $pre(v)$  to record the order in which the nodes of the tree are visited. In a similar way, by traversing  $T$  in *postorder*, each node  $v$  will get another number  $post(v)$ . These two numbers can be used to characterize the ancestor-descendant relationships as follows.

Let  $v$  and  $v'$  be two nodes of a tree  $T$ . Then,  $v'$  is a descendant of  $v$  iff  $pre(v') > pre(v)$  and  $post(v') < post(v)$ <sup>[24]</sup>.

As an example, have a look at the pairs associated with the nodes of the tree shown in Fig. 2. The first element of each pair is the preorder number of the corresponding node and the second is its postorder number. Using such labels, the ancestor-descendant relationships can be easily checked.

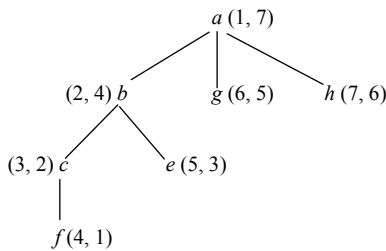


Fig. 2: Labeling a tree

For instance, by checking the label associated with  $b$  against the label for  $f$ , we see that  $b$  is an ancestor of  $f$  in terms of Proposition 1. Note that  $b$ 's label is (2, 4) and  $f$ 's label is (4, 1) and we have  $2 < 4$  and  $4 > 1$ . We also see that since the pairs associated with  $g$  and  $c$  do not satisfy the condition given in Proposition 1,  $g$  must not be an ancestor of  $c$  and *vice versa*.

Let  $(p, q)$  and  $(p', q')$  be two pairs associated with nodes  $u$  and  $v$ , respectively. We say that  $(p', q')$  is subsumed by  $(p, q)$ , denoted  $(p', q') (p, q)$ , if  $p' > p$  and

$q' < q$ . Then,  $u$  is an ancestor of  $v$  if  $(p', q')$  is subsumed by  $(p, q)$ .

In addition, if  $p' < p$  and  $q' < q$ ,  $u$  is to the left of  $v$ .

Finally, we can associate each node  $v$  with a level number  $l(v)$  (the nesting depth of the element in a document). In conjunction with the tree encoding, this number can be utilized to tell whether a node is the parent of another node. For example, if  $pre(v') > pre(v)$ ,  $post(v') < post(v)$  and  $l(v) = l(v') + 1$ , then  $v'$  is a child node of  $v$ .

### ALGORITHM FOR SIMPLE CASES

Here, we describe an algorithm for simple cases that a twig pattern contains only  $d$ -edges. First, we give a basic algorithm to show the main idea in 3.1. Then, in 3.2, we discuss how this algorithm can be substantially improved. In 3.3, we prove the correctness of the algorithm and analyze its computational complexities.

**Basic algorithm:** The basic algorithm to be given works in a bottom-up way. During the process, two data structures are maintained and computed to facilitate the discovery of subtree matchings.

- Each node  $v$  in  $T$  is associated with a set, denoted  $\alpha(v)$ , contains all those nodes  $q$  in  $Q$  such that  $Q[q]$  can be imbedded into  $T[v]$ , where  $T[v]$  represents a subtree of  $T$  rooted at  $v$ .
- Each  $q$  in  $Q$  is associated with a value  $\delta(q)$ , defined as follows:

Initially, for each  $q \in Q$ ,  $\delta(q)$  is set to  $\phi$ . During the tree matching process,  $\delta(q)$  is dynamically changed as below:

1. Let  $v$  be a node in  $T$  with parent node  $u$ .
2. If  $q$  appears in  $\alpha(v)$ , change the value of  $\delta(q)$  to  $u$ .  
Then, each time before we insert  $q$  into  $\alpha(v)$ , we will do the following checkings:
  - (i) Check whether  $label(q) = label(v)$ .
  - (ii) Let  $q_1, \dots, q_k$  be the child nodes of  $q$ . For each  $q_i$  ( $i = 1, \dots, k$ ), check whether  $\delta(q_i)$  is equal to  $v$  or to a descendant of  $v$ .

If both (1) and (2) are satisfied, insert  $q$  into  $\alpha(v)$ .

Below is a bottom-up algorithm, working in a recursive way and taking a node  $v$  in  $T$  as the input (which represents  $T[v]$ ). Initially, the input is the root of  $T$ . The algorithm will mark any node  $u$  in  $T[v]$  if it finds that  $T[u]$  contains  $Q$ . In the process, two functions are called:

- *node-check*( $u, q$ ) - It checks whether  $T[u]$  contains  $Q[q]$ . If it is the case, return  $\{q\}$ . Otherwise, it returns an empty set  $\emptyset$ .

- *leaf-node-check(u)* - It returns a set of leaf nodes in  $Q$ :  $\{q_1, \dots, q_k\}$  such that for each  $q_i$  ( $1 \leq i \leq k$ )  $label(u) = label(q_i)$ .

**Algorithm** *tree-matching(v)*

input:  $v$  - a node of tree  $T$ .

output: mark any node  $u$  in  $T[v]$  if  $T[u]$  contains  $Q$ .

**begin**

1.  $S := \emptyset; S_1 := \emptyset; S_2 := \emptyset;$
2. **if**  $v$  is not a leaf node in  $T$  **then**
3. {let  $v_1, \dots, v_k$  be the child nodes of  $v$ ;
4. **for**  $i = 1$  to  $k$  **do** call *tree-matching*( $v_i$ );
5.  $\alpha := \alpha(v_1) \cup \dots \cup \alpha(v_k);$
6. **for** each  $q \in \alpha$  **do**
7. { $\delta(q) := v; S := S \cup \{q\}$ 's parent};}
8. remove all  $\alpha(v_j)$  ( $j = 1, \dots, k$ );
9. **for** each  $q$  in  $S$  **do**
10.  $S_1 := S_1 \cup node-check(v, q);$
11. }
12.  $S_2 := leaf-node-check(v);$
13.  $\alpha(v) := \alpha \cup S_1 \cup S_2;$

**end**

**Function** *node-check(u, q)*

**begin**

1.  $S_1 := \emptyset;$
2. **if**  $label(q) = label(u)$  **then**
3. {let  $q_1, \dots, q_k$  be the child nodes of  $q$ ;
4. **if** for each  $q_i$  ( $i = 1, \dots, k$ )  $\delta(q_i)$  is equal to  $u$
5. or to a descendant of  $u$
6. **then**  $\{S_1 := S_1 \cup \{q\};$
7. **if**  $q$  is root **then** mark  $u$ ;};}
8. return  $S_1;$

**end**

**Function** *leaf-node-check(u)*

**begin**

1.  $S_2 := \emptyset;$
2. **for** each leaf node  $q$  in  $Q$  **do**
3. {**if**  $type(q) = type(u)$  **then**  $\{S_2 := \{q\};$
4. **if**  $q$  is root **then** mark  $u$ ;}
5. return  $S_2;$

**end**

The algorithm *tree-matching(v)* searches  $T$  bottom-up in a recursive way (see line 4). During the process, for each encountered node  $v$  in  $T$ , we first check whether it is a leaf node (see line 2). If it is a leaf node, the function *leaf-node-check(v)* is called (see line 12), by which all the matching leaf nodes in  $Q$  will be stored in a temporary variable  $S_2$  that will be added to  $\alpha(v)$  (see line 13). If  $v$  is an internal node, lines 3 - 10 are

first conducted and then the function *leaf-node-check(v)* is invoked (see line 12). By executing line 4, *tree-matching(v)* is recursively called for each child node  $v_i$  of  $v$ . After that, for each  $q$  appearing in  $\alpha(v_i)$ , its  $\delta$  value is set to be  $v$  (see line 7). In addition,  $q$ 's parent is inserted into  $S$ , a temporary variable to be used in a next step. Since  $\alpha(v_i)$ 's will not be used any more after this step, they are simply removed (see line 8). By executing lines 9 - 10, we check, for each  $q'$  in  $S$ , whether  $v$  matches  $q'$  by calling *node-check(v, q')*, in which the  $\delta$  values of  $q'$ 's child nodes are utilized to facilitate the checkings (see lines 3 - 5 in *node-check(v, q')*). The following example helps for illustration.

**Example 1:** Consider  $T$  and  $Q$  shown in Fig. 3.

The algorithm works in a bottom-up way. First,  $v_3$  in  $T$  is visited. It is a leaf node, matching  $q_3$  of the two leaf nodes in  $Q$ . Therefore,  $\alpha(v_3) = \{q_3\}$  (see lines 12). In the same way, we will set  $\alpha(v_5) = \{q_2\}$ .

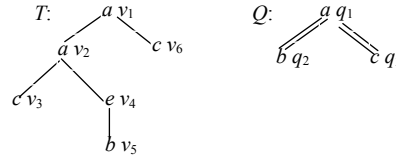


Fig. 3: A document tree and a query tree

In a next step,  $v_4$  is encountered. It is the parent of  $v_5$ . In terms of  $\alpha(v_5) = \{q_2\}$ ,  $\delta(q_2)$  is set to be  $v_4$  (Fig. 4) After that, *node-check(v4, q1)* is invoked. (Note that  $q_1$  is the parent of  $q_2$ . See lines 9 - 10.) Since  $label(v_4) \neq label(q_1)$ , it returns  $S_1 = \emptyset$ . *leaf-node-check(v4)* also returns  $S_2 = \emptyset$ . So  $\alpha(v_4) = \alpha(v_5) \cup S_1 \cup S_2 = \{q_2\}$  (see line 13). When  $v_2$  is met, we will first set  $\delta(q_2) = \delta(q_3) = v_2$  (in terms of  $\alpha(v_4) = \{q_2\}$  and  $\alpha(v_3) = \{q_3\}$ , respectively). Next, we call *node-check(v2, q1)*, in which we will check whether  $label(v_2) = label(q_1)$ . It is the case. So we will further check whether  $\delta(q_i)$  ( $i = 2, 3$ ) is equal to  $v_2$ . Since both  $\delta(q_2)$  and  $d(q_3)$  are equal to  $v_2$ , we have that  $T[v_2]$  contains  $Q[q_1]$ . Therefore,  $S_1 = \{q_1\}$ . Thus, we set  $\alpha(v_2) = \alpha(v_3) \cup \alpha(v_4) \cup S_1 \cup S_2 = \alpha(v_3) \cup \alpha(v_4) \cup \{q_1\} \cup \emptyset = \{q_1, q_2, q_3\}$

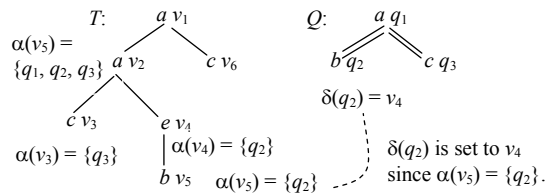


Fig. 4: Sample trace

In a next step, we will meet  $v_6$ . It is a leaf node, matching  $q_3$ . Therefore,  $\alpha(v_6) = \{q_3\}$ . Finally, we will meet  $v_1$  and set  $\delta(q_1) = v_1$  and  $\delta(q_3) = v_1$ . Since  $label(v_1) = label(q_1)$ ,  $\delta(q_2) = \delta(q_3) = v_1$ , we have that  $T[v_1]$  contains  $Q[q_1]$  and  $\alpha(v_1) = \{q_1, q_2, q_3\}$ .

**Improvements:** The above algorithm can be substantially improved by elaborating the construction of  $\alpha(v)$ 's.

First, we notice that in the case that  $v$  is a leaf node in  $T$ ,  $\alpha(v)$  is a set of the leaf nodes in  $Q$ , which match  $v$ . Such nodes can be stored in a linked list as illustrated below:

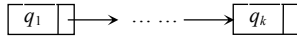


Fig. 5: Linked list to store

with the left-most node appearing first and the right-most node last. Then, for any  $1 \leq i \leq j \leq k$ , we have  $pre(q_i) < pre(q_j)$  and  $post(q_i) < post(q_j)$ . That is, in  $\alpha(v)$ ,  $q_i$ 's are sorted according to their preorder and postorder values.

Now we consider two  $\alpha$ -lists  $\alpha$  and  $\alpha'$  sorted according to their nodes' preorder and postorder numbers. Define a merging operation over  $\alpha$  and  $\alpha'$ , denoted  $merge(\alpha, \alpha')$ , as follows:

1. Assume that  $\alpha = \{v_1, \dots, v_p\}$  and  $\alpha' = \{v_1', \dots, v_q'\}$ . We step through both  $\alpha$  and  $\alpha'$  from left to right. Let  $v_i$  and  $v_j'$  be the nodes encountered. We'll make the following checkings.
2. If  $pre(v_i) > pre(v_j')$  and  $post(v_i) > post(v_j')$ , insert  $v_j'$  into  $\alpha$  after  $v_{i-1}$  and before  $v_i$  and move to  $v_{j+1}'$  (in  $\alpha'$ ).
3. If  $pre(v_i) > pre(v_j')$  and  $post(v_i) < post(v_j')$ , remove  $v_i$  from  $\alpha$  and move to  $v_{i+1}$ . (\* $v_i$  is subsumed by  $v_j'$ .\*)
4. If  $pre(v_i) < pre(v_j')$  and  $post(v_i) > post(v_j')$ , ignore  $v_j'$  and move to  $v_{j+1}'$  (in  $\alpha'$ ). (\* $v_j'$  is subsumed by  $v_i$ .\*)
5. If  $pre(v_i) < pre(v_j')$  and  $post(v_i) < post(v_j')$ , ignore  $v_i$  and move to  $v_{i+1}$ .
6. If  $pre(v_i) = pre(v_j')$  and  $post(v_i) = post(v_j')$ , ignore both  $v_i$  and  $v_j'$  and move to  $v_{i+1}$  and  $v_{j+1}'$  in  $\alpha$  and  $\alpha'$ , respectively.

The result of  $merge(\alpha, \alpha')$  is stored in  $\alpha$  and  $\alpha'$  remains unchanged. Especially, the changed  $\alpha$  is still sorted according to their nodes' preorder and postorder numbers.

In terms of the above discussion, we have the following algorithm to merge two sorted  $\alpha$ -lists together.

**Algorithm**  $merge(\alpha, \alpha')$

Input:  $\alpha$  and  $\alpha'$  - sorted  $\alpha$ -lists.

Output: modified  $\alpha$ , containing all the nodes in  $\alpha$  and  $\alpha'$  with all the subsumed nodes removed.

**begin**

1.  $p \leftarrow first\_element(\alpha)$ ;
2.  $q \leftarrow first\_element(\alpha')$ ;
3. **while**  $p \neq nil$  **do** {
4.   **while**  $q \neq nil$  **do** {
5.     **if**  $(pre(p) > pre(q) \wedge post(p) > post(q))$
6.     **then** {insert  $q$  into  $\alpha$  before  $p$ ;
7.          $q \leftarrow next(q)$ ; } (\* $next(q)$  represents the node next to  $q$  in  $\alpha'$ .\*)
8.     **else if**  $(pre(p) > pre(q) \wedge post(p) < post(q))$
9.     **then** { $p' \leftarrow p$ ; (\* $p$  is subsumed by  $q$ ; remove  $p$  from  $\alpha$ .\*)
10.         remove  $p$  from  $\alpha$ ;
11.          $p \leftarrow next(p')$ ; } (\* $next(p')$  represents the node next to  $p'$  in  $\alpha$ .\*)
12.     **else if**  $(pre(p) < pre(q) \wedge post(p) > post(q))$
13.     **then** { $q \leftarrow next(q)$ ; } (\* $q$  is subsumed by  $p$ ; move to the node next to  $q$ .\*)
14.     **else if**  $(pre(p) < pre(q) \wedge post(p) < post(q))$
15.     **then** { $p \leftarrow next(p)$ ; }
16.     **else if**  $(pre(p) = pre(q) \wedge post(p) = post(q))$
17.     **then** { $p \leftarrow next(p)$ ;
- $q \leftarrow next(q)$ ; } }
18. **if**  $p = nil \wedge q \neq nil$
- then** {attach the rest of  $\alpha'$  to the end of  $\alpha$ ; }
- end**

We can extend the merging operation over to more than two sorted  $\alpha$ -lists:

$$merge(\alpha_1, \dots, \alpha_{k-1}, \alpha_k) = merge(merge((\alpha_1, \dots, \alpha_{k-1}), \alpha_k).$$

Using this operation, the algorithm  $tree\_matching()$  is rewritten as follows.

**Algorithm**  $tree\_matching(v)$

input:  $v$  - a node of tree  $T$ .

output: mark any node  $u$  in  $T[v]$  if  $T[u]$  contains  $Q$ .

**begin**

1.  $S := \emptyset$ ;
2. **if**  $v$  is not a leaf node in  $T$  **then**
3.   {let  $v_1, \dots, v_k$  be the child nodes of  $v$ ;
4.   **for**  $i = 1$  to  $k$  **do** call  $tree\_matching(v_i)$ ;
5.    $\alpha := merge(\alpha(v_1), \dots, \alpha(v_k))$ ;
6.   assume that  $\alpha = \{q_1, \dots, q_j\}$ ;

```

7.  for  $i = 1$  to  $j$  do
8.  { $\delta(q_i) := v$ ;
9.  if ( $q_i$ 's parent  $\neq q_{i-1}$ 's parent) then
10.    $S := S \cup \{q_i$ 's parent $\}$ ;
11.  remove all  $\alpha(v_j)$  ( $j = 1, \dots, k$ );
12.  for each  $q$  in  $S$  do
13.    $S_1 := S_1 \cup \text{node-check}(v, q)$ ;
14.  }
15.  $S_2 := \text{leaf-node-check}(v)$ ;
16.  $\alpha(v) := \text{merge}(\alpha, S_1, S_2)$ ;
end

```

This algorithm is almost the same as the previous one, but with the merge operation involved, which effectively reduces the size of each  $\alpha(v)$  from  $O(|Q|)$  to  $O(Q_{leaf})$ . Special attention should also be paid to line 7, by which we generate a set  $S$  that contains the parent nodes of all those nodes appearing in  $\alpha(v_j)$ 's ( $j = 1, \dots, k$ ), where  $v_j$  is a child node of the current node  $v$ . Since the nodes in  $\alpha$  ( $\alpha = \text{merge}(\alpha_1, \dots, \alpha_{k-1}, \alpha_k)$ ) are left-to-right sorted (according to the nodes' preorder and postorder numbers), if there are more than one nodes in a sharing the same parent, they must appear consecutively in the list. So each time we insert a parent node  $q'$  (of some  $q$  in  $\alpha$ ) into  $S$ , we need to check whether it is the same as the previously inserted one. If it is the case,  $q'$  will be ignored. Thus, the size of  $S$  is also bounded by  $O(Q_{leaf})$ .

**Correctness and computational complexity:** In this subsection, we prove the correctness of the algorithm `tree-matching()` and analyze its computational complexities.

**Proposition 1:** Let  $v$  be a node in  $T$ . Then, for each  $q$  in  $\alpha(v)$  generated by `tree-matching()`, we have  $T[v]$  contains  $Q[q]$ .

**Proof:** We prove the proposition by induction on the height of  $Q$ ,  $\text{height}(Q)$ .

**Basic step:** When  $\text{height}(Q) = 1$ , the proposition trivially holds.

**Induction step:** Assume that the proposition holds for any query tree  $Q'$  with  $\text{height}(Q') \leq h$ . We consider a query tree  $Q$  of height  $h + 1$ . Let  $r_Q$  be the root of  $Q$ . Let  $q_1, \dots, q_k$  be the child nodes of  $r_Q$ . Then, we have  $\text{height}(Q[q_j]) \leq h$  ( $j = 1, \dots, k$ ). In terms of the induction hypothesis, for each  $q$  in  $Q[q_j]$  ( $j = 1, \dots, k$ ), if it appears in  $\alpha(v_i)$  (where  $v_i$  is a child node of  $v$ ), we have  $T[v_i]$  contains  $Q[q]$  and  $\delta(q)$  will be set to be  $v$ . Especially, if  $T[v_i]$  contains  $Q[q_j]$  ( $j = 1, \dots, k$ ), we have  $q_j \in \alpha(v_i)$  and

$\delta(q_j)$  will be set to be  $v$  before  $v$  is checked against  $r_Q$ . Obviously, if  $\text{label}(v) = \text{label}(r_Q)$  and for each  $q_j$  ( $j = 1, \dots, k$ ),  $\delta(q_j)$  is equal to  $v$  or a descendant of  $v$ ,  $Q$  can be embedded into  $T[v]$ . So  $r_Q$  is inserted into  $\alpha(v)$ .

Now we consider the time complexity of the algorithm, which can be divided into four parts:

1. The first part is the time spent on merging  $\alpha(v_1), \dots, \alpha(v_k)$ , where  $v_i$  ( $i = 1, \dots, k$ ) is a child node of some node  $v$  in  $T$ . This part of cost is bounded by

$$O\left(\sum_i^{|T|} d_i Q_{leaf}\right) = O(|T|Q_{leaf}),$$

where  $d_i$  represents the outdegree of a node  $v_i$  in  $T$ .

2. The second part is the time used for generating  $S$  from a merged  $\alpha$ -list. Since the size of the  $\alpha$ -list is bounded by  $O(Q_{leaf})$ , so this part of cost is also bounded by  $O(Q_{leaf})$ .

3. The third part is the time for checking a node  $v_i$  in  $T$  against each node  $q_j$  in an  $S$ . Denote  $S_i$  the set of the nodes in  $Q$ , which are checked against  $v_i$ . We estimate this part of cost by the following sum:

$$O\left(\sum_i^{|T|} \sum_j^{|S_i|} c_j\right) = O(|T|Q_{leaf}),$$

where  $c_j$  represents the outdegree of a node  $q_j$  in  $S_i$ .

4. The fourth part is the time for checking each node in  $T$  against the leaf nodes in  $Q$ . Obviously, this part of cost is bounded by

$$O\left(\sum_i^{|T|} Q_{leaf}\right) = O(|T|Q_{leaf}).$$

In terms of the above analysis, we have the following proposition.

**Proposition 2:** The time complexity of `tree-matching()` is bounded by  $O(|T|Q_{leaf})$ .

**Proof:** See the above discussion.

Since at each time point at most  $T_{leaf}$  nodes in  $T$  are associated with a  $\alpha$ -list, the space overhead is bounded by  $O(T_{leaf}Q_{leaf})$ .

## GENERAL CASES

The algorithm discussed earlier can be easily extended to general cases that a query tree contains both  $c$ -edges and  $d$ -edges. We only need to make the following changes:

- For each child node  $q_i$  of  $q$  that is being checked against  $v$ , if  $(q, q_i)$  is a  $c$ -edge, we will check whether  $\delta(q_i)$  is equal to  $v$ . If  $(q, q_i)$  is a  $d$ -edge, we simply check whether  $\text{pre}(\delta(q_i)) \geq \text{pre}(v)$  and  $\text{post}(\delta(q_i)) \leq \text{post}(v)$ .
- Accordingly, the algorithm `node-check` described earlier should be slightly modified.

**Function** *general-node-check*( $u, q$ )

**begin**

1.  $S_1 := \emptyset$ ;
2. **if**  $label(q) = label(u)$  **then**
3. {let  $q_1, \dots, q_g$  be the child nodes of  $q$ ;
4.  $flag := true; i := 1$ ;
5. **while** ( $i \leq g \wedge flag$ ) **do**
6. {**if**  $\neg(((q, q_i)$  is a  $c$ -edge)  $\wedge$  ( $\delta(q_i) = v$ ))  $\vee$
7.  $((q, q_i)$  is a  $d$ -edge)  $\wedge$
8.  $(pre(\delta(q_i)) \geq pre(u)) \wedge$
9.  $(post(\delta(q_i)) \leq post(u))$ )};
10. **then**  $flag := false$ ;
11. **if**  $i > g$  **then**  $\{S_1 := S_1 \cup \{q\}$ ;
12. **if**  $q$  is root **then** mark  $u$ ;
13. **return**  $S_1$ ;

**end**

This algorithm is similar to the function *node-check*( ). The only difference is that a general subsumption checking process is used, by which  $c$ -edges and  $d$ -edges are checked in different ways.

In addition, the lines 5 - 10 in the algorithm *tree-matching*( ) given in 3.2 should be replaced with the following segment of code:

```

for  $i = 1$  to  $k$  do {
  for  $q \in \alpha(v_i)$  do {
    let  $q'$  be the parent of  $q$ ;
    if  $((q', q)$  is a  $d$ -edge) or
       $((q', q)$  is a  $c$ -edge and  $q$  matches  $v_i$ )
    then  $\{\delta(q) := v;$ 
      let  $q''$  be the last element in  $S$ ;
      if  $(q'$ 's parent  $\neq q''$ )
      then  $S := S \cup \{q'$ 's parent};
      else remove  $q$  from  $\alpha(v_i)$ ;
    }
  }
 $\alpha := merge(\alpha(v_1), \dots, \alpha(v_k))$ ;

```

Concerning the correctness of the algorithm, we have to answer a question: whether any  $c$ -edge in  $Q$  is correctly checked.

First, we note that any  $c$ -edge in  $Q$  cannot be matched to any path with length larger than 1 in  $T$ . That is, it can be matched only to a single edge in  $T$ . It is exactly what is done by the algorithm.

Each time we check a node  $v$  in  $T$  against some  $q$  in  $Q$ , we will first set  $\delta$  values for any  $q_i$  appearing in  $\alpha(v_j)$ 's, where  $v_j$  is a child node of  $v$ . When doing this, for some  $q_i$ 's, their  $\delta$  values are changed (to  $v$ ). Assume that the current  $\delta$  value for  $q_i$  is  $v'$  (i.e.,  $\delta(q_i) = v'$ ). Then,  $v'$  must be a descendant of  $v$  since the algorithm searches  $T$  in a bottom-up way. However, we need to change  $\delta(q_i)$  from  $v'$  to  $v$  since a  $c$ -edge can match only a single edge in  $T$  and the fact that  $q_i$  matches  $v_j$  should

be recorded so that the  $c$ -edge matching is not missed (see Fig. 6 for illustration).

In Fig. 6,  $v''$  is a descendant of  $v$  and matches  $q_2$ . So  $\delta(q_2)$  will be set to  $v''$ . However,  $(q, q_2)$  is a  $c$ -edge. Therefore, the fact that  $v''$  matches  $q_2$  makes no contribution to the matching of  $v$  with  $q$ . Since  $q_2$  also matches  $v_2$ ,  $\delta(q_2)$  will be changed to  $v$ , which enables us to find that  $T[v]$  contains  $Q[q]$ .

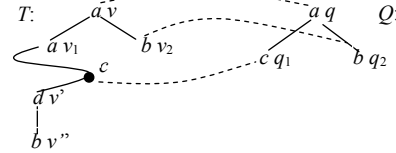


Fig. 6: Illustration for  $c$ -edge checking

In conjunction with Proposition 1, the above analysis shows the correctness of the algorithm. We have the following proposition.

**Proposition 3:** Let  $Q$  be a twig pattern containing both  $c$ -edges and  $d$ -edges. Let  $v$  be a node in  $T$ . For each  $q$  in  $\alpha(v)$  generated by *tree-matching*( ) with *general-node-check*( ) being used, we have  $T[v]$  contains  $Q[q]$ .

**Proof:** See the above discussion.

The time and space complexities for the general cases are the same as for the simple cases.

## CONCLUSION

In this article, a new algorithm is proposed for a kind of tree matching, the so-called twig pattern matching. This is a core operation for XML query processing. The main idea of the algorithm is to explore both  $T$  and  $Q$  bottom-up, by which each node  $q$  in  $Q$  is associated with a value (denoted  $\delta(q)$ ) to indicate a node  $v$  in  $T$ , which has a child node  $v'$  such that  $T[v']$  contains  $Q[q]$ . In this way, the tree embedding can be checked very efficiently. In addition, by using the tree encoding, as well as the subsumption checking mechanism, we are able to minimize the size of the lists of the matching query nodes associated with the nodes in  $T$  to reduce the space overhead. The algorithm runs in  $O(|T| \cdot Q_{leaf})$  time and  $O(T_{leaf} \cdot Q_{leaf})$  space, where  $T_{leaf}$  and  $Q_{leaf}$  represent the numbers of the leaf nodes in  $T$  and in  $Q$ , respectively. More importantly, no costly path join operation is necessary.

## REFERENCES

1. Florescu, D. and D. Kossman, 1999. Storing and querying XML data using an RDBMS. IEEE Data Eng. Bull., 22: 27-34.

2. World Wide Web Consortium. XML Path Language (XPath), W3C Recommendation, Version 1.0, Nov. 1999, <http://www.w3.org/TR/xpath>.
3. World Wide Web Consortium. XQuery 1.0: An XML Query Language, W3C Recommendation, Version 1.0, Dec. 2001, <http://www.w3.org/TR/xquery>.
4. Dutch, A., M. Fernandez, D. Florescu, A. Levy and D. Suci, 1999. A query language for XML. Proc. 8th World Wide Web Conf., pp: 77-91.
5. Chamberlin, D.D., J. Clark, D. Florescu and M. Stefanescu. XQuery1.0: An XML Query Language. <http://www.w3.org/TR/query-data-model/>.
6. Chamberlin, D.D., J. Robie and D. Florescu, 2000. Quilt: An XML query language for heterogeneous data sources. WebDB 2000.
7. Hoffmann, C.M. and M.J. O'Donnell, 1982. Pattern matching in trees. *J. ACM*, 29: 68-95.
8. Abiteboul, S., P. Buneman and D. Suci, 1999. Data on the web: From relations to semistructured data and XML. Morgan Kaufmann Publisher, Los Altos, CA 94022.
9. Chung, C., J. Min and K. Shim, 2002. APEX: An adaptive path index for XML data. *ACM SIGMOD*.
10. Cooper, B.F., N. Sample, M. Franklin, A.B. Hiltason and M. Shadmon, 2001. A fast index for semistructured data. *Proc. VLDB*, pp: 341-350.
11. Goldman, R. and J. Widom, 1997. DataGuide: Enable query formulation and optimization in semistructured databases. *Proc. VLDB*, pp: 436-445.
12. McHugh, J. and J. Widom, 1999. Query optimization for XML. *Proc. of VLDB*.
13. Li, Q. and B. Moon, 2001. Indexing and querying XML data for regular path expressions. *Proc. VLDB*, pp: 361-370.
14. Wang, H. and X. Meng, 2005. On the sequencing of tree structures for XML indexing. *Proc. Conf. Data Eng., Tokyo, Japan*, pp: 372-385.
15. Zhang, C., J. Naughton, D. Dewitt, Q. Luo and G. Lohman, 2001. On supporting containment queries in relational database management systems. *Proc. ACM SIGMOD*.
16. Kaushik, R., P. Bohannon, J. Naughton and H. Korth, 2002. Covering indexes for branching path queries. *ACM SIGMOD*.
15. Schmidt, A.R., F. Waas, M.L. Kersten, D. Florescu, I. Manolescu, M.J. Carey and R. Busse, 2001. The XML benchmark project. Technical Report INS-Ro1o3, Centrum voor Wiskunde en Informatica.
18. Bruno, N., N. Koudas and D. Srivastava, 2002. Holistic twig hoins: Optimal XML pattern matching. *Proc. SIGMOD Intl. Conf. on Management of Data, Madison, Wisconsin*, pp: 310-321.
19. Chen, T., J. Lu and T.W. Ling, 2005. On boosting holism in XML twig pattern matching. *Proc. SIGMOD*, pp: 455-466.
20. Choi, B., M. Mahoui and D. Wood, 2003. On the optimality of holistic algorithms for twig queries. *Proc. DEXA*, pp: 235-244.
21. Chen, S. *et al.*, 2006. Twig<sup>2</sup>Stack: Bottom-up processing of generalized-tree-pattern queries over XML Documents. *Proc. VLDB, Seoul, Korea*, pp: 283-323.
22. Lu, J., T.W. Ling, C.Y. Chan and T. Chan, 2005. From region encoding to extended dewey: On efficient processing of XML twig pattern matching. *Proc. VLDB*, pp: 193-204.
23. Seo, C., S. Lee and H. Kim, 2003. An Efficient Index Technique for XML Documents Using RDBMS. *Inform. Software Technol.*, 45: 11-22, Elsevier Science B.V.
24. Knuth, D.E., 1969. *The Art of Computer Programming. Vol.1*, Addison-Wesley, Reading.