

Efficient Algorithms for Descendant-Only Tree Pattern Queries[☆]

Michaela Götz^a, Christoph Koch^a, Wim Martens^{b,1}

^a*Cornell University
Ithaca, NY 14853, USA*

^b*Technical University of Dortmund
Germany*

Abstract

Tree pattern matching is a fundamental problem that has a wide range of applications in Web data management, XML processing, and selective data dissemination. In this paper we develop efficient algorithms for the tree homeomorphism problem, i.e., the problem of matching a tree pattern with exclusively transitive (descendant) edges. We first prove that deciding whether there is a tree homeomorphism is LOGSPACE-complete, improving on the current LOGCFL upper bound. Furthermore, we develop a practical algorithm for the tree homeomorphism decision problem that is both space- and time efficient. The algorithm is in LOGDCFL and space consumption is strongly bounded, while the running time is linear in the size of the data tree. This algorithm immediately generalizes to the problem of matching the tree pattern against all subtrees of the data tree, preserving the mentioned efficiency properties.

Key words: XML, XPath, query processing, tree pattern queries, complexity

1. Introduction

Tree patterns are a simple query language for tree-structured data. They are at the heart of several widely-used Web languages such as XPath and XQuery [4]. As a consequence, they form part of a number of typing mechanisms such as XML Schema, and of Web Programming Languages. They have also been used as query languages in their own right, for example for expressing subscriptions in publish-subscribe systems [1, 5, 6, 14].

[☆]The present paper is the full version of reference [13], which appeared in the Symposium on Data Base Programming Languages 2007.

Email addresses: goetz@cs.cornell.edu (Michaela Götz), koch@cs.cornell.edu (Christoph Koch), wim.martens@udo.edu (Wim Martens)

¹The third author wants to express his gratitude towards the FWO-Vlaanderen for a scholarship that permitted him to visit Christoph Koch in the Technical University of Vienna in January–February, 2005.

	time	space	streaming
Yannakakis 1981 [20]	$O(Q \cdot D \cdot \text{depth}(D))$	$O(\text{depth}(Q) \cdot D)$	no
Gottlob et al. 2002 [10]	$O(Q \cdot D)$	$O(Q \cdot D)$	no
Olteanu et al. 2004 [17]	$O(Q \cdot D \cdot \text{depth}(D))$	$O(Q \cdot \text{depth}(D) + D)$	yes
Bar-Yossef et al. 2005 [3]	$O(Q \cdot D)$	$O(Q \cdot \log D + \text{cand}_D)$	yes
Ramanan 2005 [18]	$O((Q + \text{depth}(D)) \cdot D)$	$O(Q \cdot \text{depth}(D) + \text{cand}_D)$	yes
Our bottom-up algorithm	$O(Q \cdot D \cdot \text{depth}(Q))$	$O(\text{depth}(D) \cdot \text{branch}(D))$	no
Our LOGSPACE algorithm	$\text{poly}(Q + D)$	$O(\log(Q + D))$	no

Table 1: Time and space consumption for algorithms solving the tree homeomorphism matching problem. Here $\text{depth}(\cdot)$ and $\text{branch}(\cdot)$ denote the depth and maximal branching factor of a tree, respectively.

The general tree pattern matching problem considered in the literature is the problem of finding a mapping between two node-labeled trees which is, in a sense, a cross of a subtree homomorphism and a homeomorphism. In this article we consider a clean and important special case of the tree pattern embedding problem that we call the *tree homeomorphism problem*. The question we consider is whether there is a mapping θ from the nodes of the first tree, the *tree pattern* or *query*, to the nodes of the second tree, the *data tree*, such that if node y is a child of x in the first tree, then $\theta(y)$ is a *descendant* of $\theta(x)$ in the second tree. We also consider the *tree homeomorphism matching problem*: finding *all* nodes v of the data tree such that there is such a tree homeomorphism with v the image of the root node of the pattern tree. This problem of selecting all nodes whose subtrees match the tree pattern has frequent application in XML and Web query processing [1, 10].

While this problem is of immediate practical relevance and a substantial number of papers have studied complexity and efficient algorithms for tree pattern matching, the precise complexity of both the general tree pattern matching problem and the tree homeomorphism problem are open; they are both known to be in LOGCFL and LOGSPACE-hard [11]. The former can be immediately concluded from earlier results on the complexity of the acyclic conjunctive queries [12] and the positive navigational fragment of XPath [11], both much stronger languages. The latter is a direct consequence of the fact that reachability in trees is LOGSPACE-complete [8].

Much work has been dedicated to developing efficient algorithms for finding matches of tree patterns and tree homeomorphisms. Certain algorithms aim at processing the data tree as a stream (i.e., in a single scan) [5, 6, 14, 16, 9, 17, 2, 3, 18]. For this case a number of lower bound results have been obtained using mechanisms from communication complexity [2, 3, 15]. It is basically known that streaming algorithms for even simple tree patterns consume space proportional to the size of the data tree in the worst case. Table 1 lists algorithms for the tree homeomorphism matching problem together with bounds on their running time and space consumption. Here D is the data tree and Q is the tree pattern. We assume a random access machine model with unit cost for reading and writing integers. Some of the algorithms presented support generalizations of the tree homeomorphism problem but where a better bound is known for the tree homeomorphism problem, it is shown. Some of the streaming algorithms [3, 18] use a notion of candidate node sets cand_D which depends on the algorithm and which can be of

size close to $|D|$ in the worst case. The algorithm of [3] makes the assumption of so-called non-recursive data trees, in which no two nodes such that one is a descendant of the other may have the same label. Finally, streaming algorithms such as [16] focus on being able to process SAX-events in constant time, at the cost of an exponential preprocessing step.

In this article we study the tree homeomorphism (matching) problem. We establish a tight complexity characterization and develop an algorithm for the node-selection problem (shown at the bottom of Table 1) that is both time- and space efficient. In detail, the technical contributions of this article are as follows.

- We first develop a top-down algorithm for the tree homeomorphism problem that is in LOGDCFL.²
- From this we develop a proof that the problem is LOGSPACE-complete, improving on the LOGCFL upper bound from [11].
- As our main result we develop a bottom-up LOGDCFL algorithm for computing all solutions of the tree homeomorphism problem which is both time and space efficient. This is a rather difficult algorithm and the correctness proof is involved. The algorithm runs in time $O(|D| \cdot |Q| \cdot \text{depth}(Q))$ and employs a stack of depth bounded by $\mathcal{O}(\text{depth}(D) \cdot \text{branch}(D))$.

The algorithm may be of relevance in practical implementations. Indeed, in most Web or XML applications, the data tree is *much* larger than the tree pattern yet its depth is rather small. It can be observed that ours is the only algorithm in Table 1 — and to the best of our knowledge, in existence — that can guarantee a space bound that does not contain the size, but only depth and branching factor, of the data tree as a term. At the same time the algorithm admits a good time bound.

Furthermore, the algorithm is of relevance in theory as well. It is a first step in classifying the complexity of positive Core XPath with only child and descendant axes, which is probably the most widely used XPath fragment in practice. Its precise complexity, however, is unknown.

- In some applications (e.g., for certain XML data trees), a few nodes can have a very large number of children. Our algorithm can be made to run in space $O(\text{depth}(D) \cdot \log(\text{branch}(D)))$ with the same time bound if we assume the data tree to be in a ranked form that can be obtained by a LOGSPACE linear-time preprocessing algorithm. Given that ours is an offline algorithm it means little loss of generality to assume that data trees are kept in a database in this preprocessed form.

The article presents these result basically in the order given here.

2. Definitions

By \mathbb{N} we denote the set of strictly positive integers. By Σ we always denote a fixed but infinite set of labels. The trees we consider are rooted, ordered, finite, labeled, unranked

²For our purposes, it is enough to know that LOGDCFL is characterized by deterministic logspace bounded pushdown automata which run in polynomial time [19].

trees, which are directed from the root downwards. That is, we consider trees with a finite number of nodes and in which nodes can have arbitrarily many children. A Σ -tree t (or tree t) is a relational structure over a finite number of unary labeling relations $a(\cdot)$, where each $a \in \Sigma$, and binary relations $\text{Child}(\cdot, \cdot)$ and $\text{NextSibling}(\cdot, \cdot)$. Here, $a(u)$ expresses that u is a node with label a , and $\text{Child}(u, v)$ (respectively, $\text{NextSibling}(u, v)$) expresses that v is a child (respectively, next sibling) of u . We assume that each node in a tree bears precisely one label, i.e., for each u , there is precisely one $a \in \Sigma$ such that $a(u)$ holds in t .

By ε we denote the empty tree. By $a(T_1 \cdots T_n)$ we denote the tree in which the root bears the label a and has n non-empty subtrees T_1, \dots, T_n , from left to right. If the a -labeled root has no children, we write a rather than $a()$. By $\text{root}(t)$ we denote the root node of t .

By $<_{\text{pre}}$ and $<_{\text{post}}$ we denote the depth-first left-to-right pre-ordering, respectively, left-to-right post-ordering in trees. That is, if u is a node with children u_1, \dots, u_n from left to right, then we have that $u <_{\text{pre}} u_1 <_{\text{pre}} \cdots <_{\text{pre}} u_n$ and $u_1 <_{\text{post}} \cdots <_{\text{post}} u_n <_{\text{post}} u$. Furthermore, u_1 is the successor of u in $<_{\text{pre}}$, i.e., there does not exist a v such that $u <_{\text{pre}} v <_{\text{pre}} u_1$. Similarly, u is the successor of u_n in the post-ordering. In Section 3, we will assume the $<_{\text{pre}}$ ordering on nodes, and in Section 4, we will assume the $<_{\text{post}}$ ordering.

A Σ -hedge H (or hedge H) is a finite ordered sequence $T_1 \cdots T_n$ of trees. When we write a hedge as $T_1 \cdots T_n$, we tacitly assume that every T_i is a non-empty tree. In the hedge $T_1 \cdots T_n$, we assume that $u_i <_{\text{pre}} u_{i+1}$ and $u_i <_{\text{post}} u_{i+1}$ holds for each $i = 1, \dots, n-1$, where u_i and u_{i+1} are the roots of T_i and T_{i+1} , respectively. Notice that we do not necessarily assume a sibling relation between the roots of T_i and T_{i+1} .

In the sequel, we will slightly abuse terminology and use the term “tree” to also refer to a hedge consisting of one tree, and we use the term “hedge” to also refer to the union of trees and hedges. We assume familiarity with terms such as *child*, *parent*, *descendant*, *ancestor*, *leaf*, *root*, *first child*, *last child*, *first sibling*, *previous sibling*, *last sibling*, and *next sibling*.

For a hedge H , we denote by $\text{Nodes}(H)$ the *set of nodes* of H . By $|H|$, we denote the number of nodes of H . Let $H = T_1 \cdots T_n$ with $n \geq 1$. The *label* of node u in the tree or hedge H is sometimes also denoted by $\text{lab}^H(u)$. The *depth* of a node u in H , denoted by $\text{depth}^H(u)$, is 1 when u is the root of some T_i and $1 + \text{depth}(v)$ when u is a child of v . The *height* of a node u in hedge H , denoted by $\text{height}^H(u)$, is 1 when u is a leaf and $\max(\text{height}^H(u_1), \dots, \text{height}^H(u_k)) + 1$ when u has $k > 0$ children u_1, \dots, u_k . By $\text{subtree}^H(u)$, we denote the subtree of H rooted at node u . By $\text{parent}^H(u)$, we denote the parent of u in H , if it exists. In the remainder of the article, we usually leave H implicit when H is clear from the context.

2.1. The Tree Homeomorphism Problem.

A *tree pattern query* (with descendant edges) Q is an (unranked) tree over the alphabet $\Sigma \uplus \{*\}$. That is, we assume that the special label $*$ does not appear in Σ . In the following, we use the terms *data tree* or *data hedge* to refer to ordinary Σ -trees and Σ -hedges.

Definition 1 (Tree Pattern Matching). Given a data hedge H , a node $u \in \text{Nodes}(H)$, and a tree pattern query Q , we say that H matches Q at node u , denoted by $H \models^u Q$, if there exists a mapping $h : \text{Nodes}(Q) \rightarrow \text{Nodes}(H)$ such that,

- if $\text{lab}^Q(v) = a$ for some $a \in \Sigma$, then $\text{lab}^H(h(v)) = a$;
- if $\text{Child}(v_1, v_2)$ holds in Q , then $h(v_1)$ is an ancestor of $h(v_2)$ in H ; and
- $u = h(\text{root}(Q))$.

If the above mapping h exists, we call h a *tree pattern matching*.

Notice that the ordering of children in our tree pattern queries does not matter, and that the label $*$ is a wildcard label for the query. This corresponds to the well known semantics of XPath queries with descendant axis [7]. In the following, we abbreviate by $H \models Q$ that $H \models^u Q$ for some $u \in \text{Nodes}(H)$. Alternatively, we say that H matches Q .

In this article, we are interested in the following problems.

Definition 2 (Tree Homeomorphism (Matching) Problem). Given a data tree T and a tree pattern query Q , the *tree homeomorphism problem* consists of deciding whether $T \models Q$. Furthermore, we are interested in *computing all answers* for the tree homeomorphism problem, that is, computing all nodes $u \in \text{Nodes}(T)$ such that $T \models^u Q$. We refer to the latter problem as *tree homeomorphism matching problem*.

We assume that trees are stored on tape as a set of records; one for each node. Each record contains a pointer to its first child, last child, parent, previous sibling, and next sibling.

In the remainder of the article, we assume a fixed data tree D and a fixed query tree Q for ease of presentation. We will refer to nodes of D and Q as *data nodes* and *query nodes*, respectively.

3. A Top-Down Algorithm

This section provides a simple top-down algorithm for the tree homeomorphism matching problem. The core of this top-down algorithm lies in a simple procedure that decides, given a data node d and a query node q , whether $\text{subtree}(d) \models \text{subtree}(q)$.

3.1. A Top-Down LOGDCFL Algorithm

The procedure MATCH, illustrated in Algorithm 1 tests whether $\text{subtree}(d) \models \text{subtree}(q)$. The intuition of this procedure is the following. Essentially, we immediately follow the semantics of the tree patterns. We test whether d matches q . If d matches q , it only remains to (recursively) test whether all subpatterns rooted at children of q can be matched somewhere in subtrees rooted at children d_c of the data tree d . If d does not match q , then we need to search whether $\text{subtree}(q)$ matches in some subtree rooted at some child d_c of d .

Lemma 1. MATCH is correct. That is, given a data node d and a query node q , MATCH returns true if and only if $\text{subtree}(d) \models \text{subtree}(q)$.

Algorithm 1 Top-down algorithm MATCH.

MATCH (DNode d , QNode q)
2: **if** d matches q **then**
 return \forall child q_c of $q \exists$ child d_c of d : MATCH(d_c, q_c)
4: **else** $\triangleright q$ not matched yet, try d 's children
 return \exists child d_c of d : MATCH(d_c, q)
6: **end if**

PROOF. By induction over the size of the data tree, denoted by n .

$n = 1$: We have that $\text{subtree}(d) = a$ for some $a \in \Sigma$. MATCH returns true if and only if the query tree consists of one node and d matches this node. The correctness follows from the tree pattern matching definition, which says that if $\text{subtree}(d) = a$, $\text{subtree}(q) = a$ or $\text{subtree}(q) = *$, $\text{subtree}(d) \models \text{subtree}(q)$.

$n > 1$: We consider two cases.

- If d matches q , we return true if, for every child q_c of q , there exists a child d_c of d such that MATCH(d_c, q_c) returns true. If the query tree consists of only one node, this is obviously correct. If q has children, the correctness follows from the induction hypothesis and the definition of tree pattern matchings: If $\text{subtree}(d) = a(T_1 \cdots T_n)$, $\text{subtree}(q) = x(Q_1 \cdots Q_m)$, $x \in \Sigma \uplus \{*\}$, $a \models x$, and, for every $k = 1, \dots, m$, there exists an $i_k \in \{1, \dots, n\}$, such that $T_{i_k} \models Q_k$, then $\text{subtree}(d) \models \text{subtree}(q)$. If there exists a q_c such that MATCH(d_c, q_c) is false for every d_c , we would also fail to match the whole query tree into a subtree of a child of d . Again by the definition of tree pattern matchings it is then correct to return false.
- If d does not match q , we test whether there is a child d_c of d such that $\text{subtree}(q)$ can be matched into $\text{subtree}(d_c)$. By the induction hypothesis, the recursive calls of MATCH(d_c, q) compute this correctly. If there is such a matching, it is correct to return true by the definition of tree pattern matchings: If $\text{subtree}(d) = a(T_1 \cdots T_n)$ and $T_i \models \text{subtree}(q)$, then $\text{subtree}(d) \models \text{subtree}(q)$. Furthermore, if $\text{subtree}(d) = a(T_1 \cdots T_n)$, d does not match q , and there does not exist a T_i such that $T_i \models \text{subtree}(q)$, then, by definition, $\text{subtree}(d) \not\models \text{subtree}(q)$. Hence, it is correct to return false. \square

Hence, MATCH is a correct algorithm for the tree homeomorphism problem. By slightly adapting MATCH, we can even turn it into an algorithm TOP-DOWN-MATCH for the tree homeomorphism *matching* problem too. First, we need a procedure EXACT-MATCH that, given a data node d and query node q , decides whether $\text{subtree}(d)$ matches $\text{subtree}(q)$ at node d . This is easy: EXACT-MATCH only differs from MATCH in line 5, where it just returns false. Given a data node d and the root q_{root} of the query tree, TOP-DOWN-MATCH now simply iterates over all the data nodes and returns every data node d for which EXACT-MATCH(d, q_{root}) returns true. From this construction and from the correctness of MATCH, it is now immediate that TOP-DOWN-MATCH is correct as well.

Time and Space Complexity.

We start with an analysis of the time complexity of MATCH and then we describe how an upper bound of the runtime of EXACT-MATCH can be derived from that.

Observation 1. MATCH(d, q) compares each node in subtree(d) at most once with each node in subtree(q). The running time of MATCH(d, q) is $|\text{subtree}(d)| \cdot |\text{subtree}(q)|$.

PROOF. This is an easy induction on $|\text{subtree}(d)|$. If $|\text{subtree}(d)| = 1$, then MATCH tests whether d matches q and discovers that there are no children of d to iterate over. Hence, the running time is in $\mathcal{O}(|\text{subtree}(q)|)$.

If $|\text{subtree}(d)| > 1$, then MATCH tests whether d matches q and it either calls itself recursively for every child d_c of d and every child q_c of q ; or it calls itself recursively for every child d_c of d and q . In both cases, we can apply the induction hypothesis. In the first case, the time complexity becomes $\mathcal{O}(\sum_{q_c} (\sum_{d_c} (|\text{subtree}(d_c)| \cdot |\text{subtree}(q_c)|)))$, and in the second case, the time complexity becomes $\mathcal{O}(\sum_{d_c} (|\text{subtree}(d_c)| \cdot |\text{subtree}(q)|))$. Hence, both cases are in $\mathcal{O}(|\text{subtree}(d)| \cdot |\text{subtree}(q)|)$. \square

It is easy to see that Observation 1 implies that the time complexity of EXACT-MATCH(d, q) is also in $\mathcal{O}(|\text{subtree}(d)| \cdot |\text{subtree}(q)|)$. As TOP-DOWN-MATCH simply calls EXACT-MATCH for every data node, we immediately have the following result.

Proposition 1. The running time of TOP-DOWN-MATCH is in $\mathcal{O}(|D|^2 \cdot |Q|)$. Moreover, TOP-DOWN-MATCH makes $\mathcal{O}(|D|^2 \cdot |Q|)$ comparisons between a data node and a query node.

It is immediate from our implementation that the algorithm can be executed by a deterministic logarithmic space bounded auxiliary pushdown automaton (see, e.g., [19]). Moreover, by Proposition 1, this auxiliary pushdown automaton runs in polynomial time. It follows from [19] that the tree homeomorphism matching problem is in LOGDCFL. As the maximum recursion depth of Algorithm 1 is $\mathcal{O}(\text{depth}(D))$, this renders the algorithm quite space-efficient, but the running time being quadratic in the size of the data tree, and the many unnecessary comparisons between query and data nodes are quite unsatisfactory. In Section 4, we show how these issues can be resolved by turning to a bottom-up approach.

3.2. A LOGSPACE Procedure

While the top-down algorithm does not seem to be well-suited for efficiently computing *all* nodes u for which $D \models^u Q$, it is quite useful for *deciding* whether $D \models Q$, from a complexity theory point of view. Indeed, as we will exhibit, a modified version of MATCH can decide in LOGSPACE whether $D \models Q$.

For ease of presentation of the algorithm, we assume the *depth-first left-to-right pre-order* ordering on nodes in trees and hedges in the remainder of this section. For a node u , we denote by $u + 1$ the successor node of u in the left-to-right pre-order $<_{\text{pre}}$. We note that this assumption does not restrict our algorithm as one can compute this successor in LOGSPACE.

We argue how to transform Algorithm 1 into a LOGSPACE algorithm that decides whether $D \models Q$. We will first give an intuition of the transformation. Then we will

Algorithm 2 Top-down algorithm L-MATCH. Here, $+1$ denotes the successor in the depth-first left-to-right pre-ordering.

```

L-MATCH (DNode  $d$ , QNode  $q$ )
2: if  $d$  matches  $q$ , and both  $d$  and  $q$  have children then                                 $\triangleright \theta(q) = d$ 
    return L-MATCH ( $d + 1, q + 1$ )
4: else if  $d$  does not match  $q$  and  $d$  has a child then
    return L-MATCH ( $d + 1, q$ )
6: else if  $d$  matches  $q$  and  $q$  is a leaf then                                         $\triangleright \theta(q) = d$ 
    if  $q$  is maximal then
8:     return true                                                                     $\triangleright$  none of  $q$ 's ancestors has a next sibling
    else
10:     $d' \leftarrow$  BACKTRACK( $d, q + 1$ )                                              $\triangleright$  node to which parent( $q + 1$ ) matched
        return L-MATCH ( $d' + 1, q + 1$ )
12:    end if
    else                                                                               $\triangleright d$  is a leaf and ( $d$  does not match  $q$  or  $q$  is not a leaf)
14:    if  $d$  is maximal then
        return false
16:    end if
         $q' \leftarrow q$ 
18:    while  $q'$  has a parent do
         $d' \leftarrow$  BACKTRACK( $d, q'$ )                                              $\triangleright$  node to which parent( $q'$ ) was matched
20:        if  $d'$  is an ancestor of  $d + 1$  then
            return L-MATCH ( $d + 1, q'$ )
22:        else  $q' \leftarrow$  parent( $q'$ )
            end if
24:    end while
        return L-MATCH ( $d + 1, q'$ )
26: end if

```

discuss some implementation details that will allow us to analyze the space consumption. A formal proof of the correctness follows.

Intuitively, the LOGSPACE algorithm processes the data and query trees in a top-down manner, just like Algorithm 1, and it processes the children of a node from left to right. Whenever Algorithm 1 uses the recursion stack to determine which function call to issue next or which final value to return, the LOGSPACE algorithm *recomputes* the information necessary to make these decisions.

Therefore, the essential difference between Algorithm 1 and the LOGSPACE algorithm lies in a backtracking procedure. When, for example, Algorithm 1 matches a leaf q of the query tree onto some data node d , then it uses the recursion stack to discover the data node onto which q 's parent was matched in the data tree and tries to match q 's next sibling in some subtree of that data node. Instead of using this recursion stack, the LOGSPACE algorithm enters a subprocedure BACKTRACK(d, q) that *recomputes* the data node onto which q 's parent was matched. In particular, BACKTRACK(d, q) computes the highest possible node d' on the path from D 's root to d , such that the path from D 's root to d' matches the path from Q 's root to q 's parent. The crux of the algorithm is that this is *correct*, i.e., d' is equal to the data node onto which q 's parent was matched;

and that $\text{BACKTRACK}(d, q)$ can be performed using only logarithmic space on a Turing Machine. $\text{BACKTRACK}(d, q)$ stores d and q on tape and goes to the roots of the query and data tree. It then matches the path to d with the path to q in a greedy manner. The crux of executing $\text{BACKTRACK}(d, q)$ using logarithmic space lies in the following. If we arrive at a node u in D (resp., Q), we have to be able to determine the child of u that lies on the path to d (resp., q). To this end, we first store d (resp., q) in a temporary variable v . We continue following the parent relation in this fashion until we find u , at which point we return the value of v , which is a child of u .³

In more detail, for given input nodes d and q the LOGSPACE procedure tests whether d matches q and based on the result of this test it computes the next function call. This is a rather extensive case study. In case d matches q and both nodes have children the next function call has the leftmost child of d and the leftmost child of q as its input. In case d does not match q but has children the next function call has the leftmost child of d and q as its input. In other cases, computing the next function call can be more complicated. When, for example, Algorithm 1 matches a leaf q of the query tree onto some data node d it will try to match $q + 1$ next, which is the lowest right sibling we encounter on the path from q to the root. If no such sibling exists, all query nodes are matched and the algorithm returns true. Otherwise, Algorithm 1 uses the recursion stack to compute the data node onto which $q + 1$'s parent was matched in the data tree and tries to match $q + 1$ in some proper subtree of that data node. Instead of using this recursion stack, the LOGSPACE algorithm enters the subprocedure $\text{BACKTRACK}(d, q + 1)$ that recomputes the data node onto which $q + 1$'s parent was matched. The next function call in that case has the leftmost child of $\text{BACKTRACK}(d, q + 1)$ and $q + 1$ as input. There is one more case: If d is a leaf and either d does not match q or q has children, then Algorithm 1 tries to match q to d 's right sibling if it has one. In general, Algorithm 1 will try to move a query node onto $d + 1$ next if such a node exists, otherwise it returns false. If $d + 1$ exists, it uses the recursion stack to find the ancestor-or-self of q that is closest to q and whose parent was matched to an ancestor of $d + 1$. Algorithm 1 tries to match this ancestor in $\text{subtree}(d + 1)$. If no such parent exists then Algorithm 1 tries to match the root in the $\text{subtree}(d + 1)$. Analogously as before, the LOGSPACE algorithm uses BACKTRACK to test for an ancestor of q whether its parent was matched to an ancestor of $d + 1$.

We present the LOGSPACE procedure in Algorithm 2. For ease of presentation, we have written the algorithm as a recursive procedure, but it can be implemented to only use logarithmic space. This can be seen by observing that every recursive call to L-MATCH in Algorithm 2 is a return-statement, so the algorithm does not change when the recursion stack is not used at all. The input of the algorithm is, just as before, the root nodes d and q of the data tree D and query tree Q , respectively. In particular, we can rewrite the LOGSPACE procedure into a non-recursive algorithm: We wrap a while loop (with condition true) around the function body. In the function body we replace each function call by an update of d and q (according to the input of the function call) followed by a break statement. Thus we start an execution of the while loop for each function call.

For the sake of understanding the general idea behind Algorithm 2, let, for a query node q , the *remainder of q in Q* be the subhedge of Q consisting of the nodes $\{q' \mid q \leq_{\text{pre}} q'\}$

³Notice that the parent pointer is not mandatory for this argument. One can also determine v 's parent in LOGSPACE by scanning the input tape and searching for a node with a child-pointer to v .

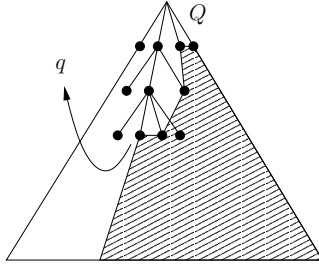


Figure 1: Illustration of the remainder of q in Q .

$q' \leq_{\text{pre}} q_{\text{max}}\}$, where q_{max} is the maximal query node w.r.t. the depth-first left-to-right ordering. We illustrate the remainder of q in Q in Figure 1. Given a data node d and a query node q , the algorithm first tries to match the remainder of q in Q consistently with what has already been matched in D (lines 2–12). If this fails, it either returns false (line 15), or enters the backtracking procedure (lines 18–25).

We argued above that we can implement BACKTRACK in LOGSPACE. Algorithm 2 does not require a recursion stack and only uses logarithmic space. Thus we have the following proposition.

Proposition 2. *Algorithm 2 runs in LOGSPACE.*

3.2.1. Correctness of L-MATCH

We want to show that L-MATCH returns true on input D and Q if and only if $D \models Q$. To simplify the analysis, we imaginarily extend the algorithm by defining a matching θ . If the algorithm compares the labels of d and q in the function call L-MATCH(d, q) and they agree (in lines 2 and 6), we set $\theta(q) = d$ (and may overwrite older assignments). This mapping θ is merely used to simplify the reasoning about the algorithm.

Soundness. We will prove that whenever L-MATCH returns true on input D and Q , then $D \models Q$. In fact we prove a stronger claim: If L-MATCH returns true, then our mapping θ is a tree pattern matching (cfr. Definition 1). Hence θ witnesses that if L-MATCH returns true, then $D \models Q$.

In order to prove the soundness of L-MATCH, we first show the following Lemma, that also implies that BACKTRACK is indeed correct. That is, given q and d , the node onto which q 's parent was matched can be computed by calculating the highest possible node d' on the path from D 's root to d , such that the path from D 's root to d' matches the path from Q 's root to q 's parent.

Lemma 2. *Let D be a data tree and Q be a query tree. Further, let L-MATCH(d, q) be a function call resulting from the initial procedure call L-MATCH($\text{root}(D), \text{root}(Q)$). Then at the time when L-MATCH(d, q) is called*

- (1) *the restriction of θ to query nodes smaller than q in the ordering $<_{\text{pre}}$ is a tree pattern matching;*
- (2) *θ matches the path $\langle \text{parent}(q) \cdots \text{root}(Q) \rangle$ into the path $\langle \text{parent}(d) \cdots \text{root}(D) \rangle$ as high as possible; and*

(3) the path $\langle q \cdots \text{root}(Q) \rangle$ cannot be matched into the path $\langle \text{parent}(d) \cdots \text{root}(D) \rangle$.

PROOF. We prove the Lemma by induction on the position k of $\text{L-MATCH}(d, q)$ in the sequence of function calls resulting from the initial procedure call $\text{L-MATCH}(\text{root}(D), \text{root}(Q))$. If $k = 1$ then we have $\text{L-MATCH}(\text{root}(D), \text{root}(Q))$, in which case there is nothing to show.

So, from now on, we assume that the Lemma is true for the first k function calls and we let $\text{L-MATCH}(d, q)$ be the k^{th} function call. We prove that it is also true for the $k + 1^{\text{th}}$ function call (if there is one). We consider four cases according to Algorithm 2.

- If the labels of d and q agree and both nodes have children (line 2), the next function call is $\text{L-MATCH}(d + 1, q + 1)$, where $d + 1$ and $q + 1$ are the leftmost children of d and q , respectively. We know by induction that θ , restricted to query nodes smaller than q , is a tree pattern matching. We extend this mapping by $\theta(q) = d$. This mapping clearly preserves labels. Hence, we only need to show that $\theta(q)$ is a descendant of $\theta(\text{parent}(q))$. But this is clear, since by induction $\langle \text{parent}(q) \cdots \text{root}(Q) \rangle$ is matched as high as possible into the path $\langle \text{parent}(d) \cdots \text{root}(D) \rangle$, which proves (1). Combining this with the fact that $\langle q \cdots \text{root}(Q) \rangle$ cannot be matched into $\langle \text{parent}(d) \cdots \text{root}(D) \rangle$ we conclude that $\langle q \cdots \text{root}(Q) \rangle$ is matched as high as possible into $\langle d \cdots \text{root}(D) \rangle$, which proves (2). As we must match $q + 1$ onto a descendant of $\theta(q)$, it then follows that the path $\langle d \cdots \text{root}(D) \rangle$ cannot match the path $\langle q + 1 \cdots \text{root}(Q) \rangle$, which proves (3).
- If the labels of d and q do not agree and d has children (line 4), the next function call is $\text{L-MATCH}(d + 1, q)$, where $d + 1$ is the leftmost child of d . We do not extend θ in that case and all requirements (1), (2), and (3) follow from the induction hypothesis.
- If the labels of d and q agree, q is a leaf, and q is not maximal (line 6), we extend the mapping θ by $\theta(q) = d$. As in the first case of this proof we know by induction that θ , restricted to query nodes less than q , is a tree pattern matching. The extended θ is still a tree pattern matching, because, due to the induction hypothesis, $\langle \text{parent}(q) \cdots \text{root}(Q) \rangle$ is matched into the path $\langle \text{parent}(d) \cdots \text{root}(D) \rangle$. Hence, (1) is true.

$\text{BACKTRACK}(d, q + 1)$ calculates the highest ancestor d' of the data node d such that $\langle d' \cdots \text{root}(D) \rangle$ matches $\langle \text{parent}(q + 1) \cdots \text{root}(Q) \rangle$. Why does d' exist? First, note that $\text{parent}(q + 1)$ is an ancestor of q due to the left-to-right pre-order ordering. Second, by induction, $\langle \text{parent}(q) \cdots \text{root}(Q) \rangle$ can be matched into the path $\langle \text{parent}(d) \cdots \text{root}(D) \rangle$. Putting both facts together, the sub-path $\langle \text{parent}(q + 1) \cdots \text{root}(Q) \rangle$ can still be matched into the path $\langle \text{parent}(d) \cdots \text{root}(D) \rangle$. Hence, d' exists and $d' + 1$ is its leftmost child.

The next function call is $\text{L-MATCH}(d' + 1, q + 1)$. By induction, the mapping θ matches the path $\langle \text{parent}(q) \cdots \text{root}(Q) \rangle$ into the path $\langle \text{parent}(d) \cdots \text{root}(D) \rangle$ as high as possible and therefore, θ also matches the sub-path $\langle \text{parent}(q + 1) \cdots \text{root}(Q) \rangle$ as high as possible into $\langle \text{parent}(d) \cdots \text{root}(D) \rangle$. It also follows that BACKTRACK in fact calculated the node onto which $\text{parent}(q + 1)$ was matched, e.g. $d' = \theta(\text{parent}(q + 1))$. Combining the last two facts with the descendant requirement that is fulfilled by θ yields (2) and (3): θ matches the sub-path $\langle \text{parent}(q + 1) \cdots \text{root}(Q) \rangle$ as high as possible into the path $\langle \text{parent}(d) \cdots \text{root}(D) \rangle$.

$\dots\text{root}(Q)\rangle$ as high as possible into the path $\langle d' \dots \text{root}(D)\rangle$ and therefore $\langle q + 1 \dots \text{root}(Q)\rangle$ cannot be matched into $\langle d' \dots \text{root}(D)\rangle$.

- If d is a leaf and (d does not match q or q is not a leaf) and d is not maximal (line 13), we have to try to match q somewhere else. We do not extend θ , so θ restricted to query nodes smaller than q is still a tree pattern matching, which proves (1). To prove the other items, we consider two cases.

(Case 1): Assume that the next function call is $\text{L-MATCH}(d + 1, q')$ in line 25. Then q' has no parent ($q' = \text{root}(Q)$) and (2) is trivially true. To prove (3), i.e., to prove that $\text{root}(Q)$ cannot be matched into the path $\langle \text{parent}(d + 1) \dots \text{root}(D)\rangle$, we consider two cases.

- If $q = q' = \text{root}(Q)$, by induction, $\text{root}(Q)$ cannot be matched into $\langle \text{parent}(d) \dots \text{root}(D)\rangle$ and therefore also not into $\langle \text{parent}(d + 1) \dots \text{root}(D)\rangle$, which is a sub-path of $\langle \text{parent}(d) \dots \text{root}(D)\rangle$, which proves (3).
- If $q \neq q' = \text{root}(Q)$, then $\text{root}(Q)$ is an ancestor of q . By the induction hypothesis on (1) we have that $\langle \text{parent}(\theta(\text{root}(Q))) \dots \text{root}(D)\rangle$ is a sub-path of $\langle d \dots \text{root}(D)\rangle$. Also, $\langle \text{parent}(d + 1) \dots \text{root}(D)\rangle$ is a sub-path of $\langle d \dots \text{root}(D)\rangle$. As L-MATCH did not return a function call in line 21, $\theta(\text{root}(Q))$ is not an ancestor of $\text{parent}(d + 1)$. Hence, $\langle \text{parent}(\theta(\text{root}(Q))) \dots \text{root}(D)\rangle$ includes $\langle \text{parent}(d + 1) \dots \text{root}(D)\rangle$. By induction, $\langle \text{parent}(\theta(\text{root}(Q))) \dots \text{root}(D)\rangle$ does not match $\text{root}(Q)$ and this property carries over to $\langle \text{parent}(d + 1) \dots \text{root}(D)\rangle$, which proves (3).

(Case 2): Otherwise, the next function call is $\text{L-MATCH}(d + 1, q')$ in line 21. $\text{BACKTRACK}(d, q')$ has calculated the highest ancestor d' of the data node d such that $\langle d' \dots \text{root}(D)\rangle$ matches $\langle \text{parent}(q') \dots \text{root}(Q)\rangle$. Why does d' exist? First, note that q' lies on the path $\langle q \dots \text{root}(Q)\rangle$ and has a parent (line 20). Further, note that, by induction, the path $\langle \text{parent}(d) \dots \text{root}(D)\rangle$ matches the path $\langle \text{parent}(q) \dots \text{root}(Q)\rangle$ and therefore it also matches the sub-path $\langle \text{parent}(q') \dots \text{root}(Q)\rangle$. It follows that d' exists and that $d' + 1$ is its leftmost child.

We know that q' is the lowest node on $\langle q \dots \text{root}(Q)\rangle$ such that $\text{BACKTRACK}(d, q') = d'$ is an ancestor of $d + 1$, by the condition in the while loop. Next, we will prove (2). By induction, the mapping θ matches the query path $\langle \text{parent}(q) \dots \text{root}(Q)\rangle$ and therefore also the sub-path $\langle \text{parent}(q') \dots \text{root}(Q)\rangle$ as high as possible into the data path $\langle \text{parent}(d) \dots \text{root}(D)\rangle$. It follows that the mapping θ also matches the path $\langle \text{parent}(q') \dots \text{root}(Q)\rangle$ as high as possible into the sub-path $\langle d' \dots \text{root}(D)\rangle$. As d' is an ancestor of $d + 1$ (line 21) we now have that the mapping θ matches the path $\langle \text{parent}(q') \dots \text{root}(Q)\rangle$ as high as possible into the path $\langle \text{parent}(d + 1) \dots \text{root}(D)\rangle$, which proves (2).

In order to prove (3), i.e., to prove that the path $\langle \text{parent}(d + 1) \dots \text{root}(D)\rangle$ cannot match the path $\langle q' \dots \text{root}(Q)\rangle$, we consider two cases:

- If $q = q'$, by the induction hypothesis, the path $\langle \text{parent}(d) \dots \text{root}(D)\rangle$ cannot match the path $\langle q \dots \text{root}(Q)\rangle$. We have that $\langle \text{parent}(d + 1) \dots \text{root}(D)\rangle$ is a sub-path of $\langle \text{parent}(d) \dots \text{root}(D)\rangle$ because d is a leaf. The claim follows.

- If $q \neq q'$, recall that q' is the lowest ancestor of q such that $\theta(\text{parent}(q'))$ is an ancestor of $d + 1$ (observe the while loop and recall that, by induction, $d' = \theta(\text{parent}(q'))$). It follows, that q' is matched somewhere on the path from $\text{parent}(d)$ to (but not including) $\text{parent}(d+1)$. By the induction hypothesis, we cannot match the path $\langle q' \cdots \text{root}(Q) \rangle$ any higher. Hence, the path $\langle \text{parent}(d+1) \cdots \text{root}(D) \rangle$ does not match the path $\langle q' \cdots \text{root}(Q) \rangle$.

• Otherwise there does not follow a function call. \square

Proposition 3. *Algorithm 2 is sound. That is, given a data D and query tree Q , if Algorithm 2 returns true, then $D \models Q$.*

PROOF. If $\text{L-MATCH}(d, q)$ returns true in line 8, then q is maximal (line 7) and the label of d matches the one of q (line 6). By Lemma 2 the mapping θ is a tree pattern matching of $Q \setminus \{q\}$ on D , such that q 's parent is matched onto some ancestor of d . We extend the mapping by $\theta(q) = d$, and conclude that $D \models Q$. \square

Completeness. In this section we want to prove that whenever L-MATCH returns false on input D and Q , then $D \not\models Q$. For two nodes x and y in a tree, we denote by $\langle x \cdots \bar{y} \rangle$ the path from x to y that excludes y itself. In order to prove the completeness, we first show the following Lemma. Recall that the previous sibling of a node is its sibling to the left.

Lemma 3. *Let D be a data tree and let Q be a query tree. Let $\text{L-MATCH}(d, q)$ be a function call resulting from the initial procedure call $\text{L-MATCH}(\text{root}(D), \text{root}(Q))$. Then, it holds for all previous siblings \hat{d} of nodes on the path $\langle d \cdots \bar{\theta(\text{parent}(q))} \rangle$ or, in case q has no parent, on the path $\langle d \cdots \text{root}(D) \rangle$ that*

$$\text{subtree}(\hat{d}) \not\models \text{subtree}(q).$$

PROOF. Note that, by Lemma 2, we can refer to the restriction of θ to query nodes smaller than q as a tree pattern matching. The proof is by induction on the position k of $\text{L-MATCH}(d, q)$ in the sequence of function calls resulting from the initial procedure call $\text{L-MATCH}(\text{root}(D), \text{root}(Q))$. If $k = 1$ then we have $\text{L-MATCH}(\text{root}(D), \text{root}(Q))$ in which case there is nothing to show because there are no left siblings on the path $\langle \text{root}(D) \rangle$.

So, from now on, we assume that $k \geq 1$ and the Lemma is true for the first k function calls. Let $\text{L-MATCH}(d, q)$ be the k^{th} function call. We prove that it is also true for the $k + 1^{\text{th}}$ function call (if there is one). We consider four cases according to Algorithm 2.

- If the labels of d and q agree and both nodes have children (line 2), $\theta(q)$ is defined to be d . The next function call is $\text{L-MATCH}(d + 1, q + 1)$, where $d + 1$ and $q + 1$ are the leftmost children of d and q , respectively. The path $\langle d + 1 \cdots \bar{\theta(\text{parent}(q + 1))} \rangle$ is the path $\langle d + 1 \cdots \bar{d} \rangle$. Since $d + 1$ has no left sibling there is nothing to show.
- If the labels of d and q do not agree and d has children (line 4), the next function call is $\text{L-MATCH}(d + 1, q)$, where $d + 1$ is the leftmost child of d . Since $d + 1$ has no left siblings, the claim follows from the induction hypothesis.

- If the labels of d and q agree, q is a leaf (line 6), and q is not maximal, $\theta(q)$ is defined to be d . $\text{BACKTRACK}(d, q + 1)$ calculates the highest ancestor d' of d such that $\langle d' \cdots \text{root}(D) \rangle$ matches $\langle \text{parent}(q + 1) \cdots \text{root}(Q) \rangle$. By Lemma 2 we have that $\theta(\text{parent}(q + 1)) = d'$. The next function call is $\text{L-MATCH}(d' + 1, q + 1)$. The path $\langle d' + 1 \cdots \theta(\text{parent}(q + 1)) \rangle$ is the path $\langle d' + 1 \cdots \hat{d}' \rangle$, where d' is $d' + 1$'s parent. As $d' + 1$ has no left sibling, there is nothing to show.
- If d is a leaf and (the labels of d and q do not agree or q has children) (line 13) and d is not maximal, then $\text{subtree}(d)$ does not match $\text{subtree}(q)$. We first show the following invariant which we will need later:

Invariant 2. *For every call of L-MATCH until the k -th call, whenever the body of the while loop in line 18 is executed without returning a function call in line 21, it follows for the current q' that $\text{subtree}(\theta(\text{parent}(q')))$ does not match $\text{subtree}(\text{parent}(q'))$.*

PROOF. We prove the claim by induction over the number of executions of the while body, denoted by ℓ .

$\ell = 1$: Here $q' = q$, q has a parent (line 18), and we know that (i) $\text{subtree}(q)$ cannot be matched into $\text{subtree}(d)$ (line 13), (ii) q cannot be matched into the path $\langle \text{parent}(d) \cdots \theta(\text{parent}(q)) \rangle$ by Lemma 2, (iii) there are no right siblings on the path $\langle d \cdots \theta(\text{parent}(q)) \rangle$, since otherwise we would have returned a function call in line 21, and (iv) $\text{subtree}(q)$ cannot be matched into $\text{subtree}(\hat{d})$ for every left sibling \hat{d} of the path $\langle d \cdots \theta(\text{parent}(q)) \rangle$, by the induction hypothesis of Lemma 3. From (i-iv) we can conclude that no proper subtree of $\theta(\text{parent}(q))$ matches $\text{subtree}(q)$, which implies that $\text{subtree}(\theta(\text{parent}(q)))$ does not match $\text{subtree}(\text{parent}(q))$.

$\ell > 1$: Let the claim be true for the first ℓ while loop executions. We prove that it is also true for the $\ell + 1^{\text{th}}$ execution. Let q' be the query node of the $\ell + 1^{\text{th}}$ while loop execution. Here, $q' \neq q$ and q' has a parent (line 18). There must have been a function call $\text{L-MATCH}(q', \theta(q'))$ and there must have been a while loop execution with the child of q' on the path from q to q' as current node. We know that (i) $\text{subtree}(q')$ cannot be matched into $\text{subtree}(\theta(q'))$ by the induction hypothesis, (ii) q' cannot be matched into the path $\langle \text{parent}(\theta(q')) \cdots \theta(\text{parent}(q')) \rangle$ by Lemma 2, (iii) there are no right siblings on the path $\langle \theta(q') \cdots \theta(\text{parent}(q')) \rangle$, since otherwise we would have returned a function call in line 21, and (iv) $\text{subtree}(q')$ cannot be matched into $\text{subtree}(\hat{d})$, for every left sibling \hat{d} of the path $\langle \theta(q') \cdots \theta(\text{parent}(q')) \rangle$ by the induction hypothesis of Lemma 3. From (i-iv), we can conclude that no proper subtree of $\theta(\text{parent}(q'))$ matches $\text{subtree}(q')$, which implies that the $\text{subtree}(\theta(\text{parent}(q')))$ does not match the $\text{subtree}(\text{parent}(q'))$. \square

We return to the proof of the main induction. We denote the left sibling of $d + 1$ by $\text{prevSib}(d + 1)$. We consider two cases.

(Case 1): Assume that next function call is $\text{L-MATCH}(d + 1, q')$ in line 25. Here, q' is the query root. We need to show that there is no left sibling \hat{d} on the path $\langle d + 1 \cdots \text{root}(D) \rangle$, such that $\text{subtree}(\hat{d}) \models \text{subtree}(q')$. We consider two cases:

- If $q = q' = \text{root}(Q)$, by the induction hypothesis, $\text{subtree}(q)$ cannot be matched into $\text{subtree}(\hat{d})$ for any left sibling \hat{d} of the path $\langle d \cdots \text{root}(D) \rangle$. Since $\text{parent}(d+1)$ is an ancestor of d , it is enough to show that the subtree rooted at $\text{prevSib}(d+1)$, which is a subtree that includes d , does not match $\text{subtree}(q)$. We know that (i) $\text{subtree}(q)$ cannot be matched into $\text{subtree}(d)$, (ii) q cannot be matched into the path $\langle \text{parent}(d) \cdots \text{root}(D) \rangle$ by Lemma 2, (iii) there are no right siblings on the path $\langle d \cdots \text{prevSib}(d+1) \rangle$ due to the left-to-right pre-order successor, and (iv) $\text{subtree}(q)$ cannot be matched into the subtree $\langle \hat{d} \rangle$ for every left sibling \hat{d} of the path from d to $\text{root}(D)$ by the induction hypothesis. From (i–iv) it follows that we cannot match $\text{subtree}(q)$ into the subtree rooted at $\text{prevSib}(d+1)$.
- If $q \neq q' = \text{root}(Q)$, then by Lemma 2 there must have been a function call $\text{L-MATCH}(q', \theta(q'))$. By the induction hypothesis, $\text{subtree}(q')$ cannot be matched into $\text{subtree}(\hat{d})$ for any of the left siblings \hat{d} of the path $\langle \theta(q') \cdots \text{root}(D) \rangle$. Furthermore, there must have been a while loop execution with q' 's child on the path from q to q' as current query node. Since $\text{parent}(d+1)$ is an ancestor of $\theta(q')$ (otherwise we would have returned a function call in line 21), it is enough to show that the subtree rooted at $\text{prevSib}(d+1)$ does not match the subtree $\langle q' \rangle$. We know that (i) $\text{subtree}(q')$ cannot be matched into $\text{subtree}(\theta(q'))$ by Invariant 2, (ii) q' cannot be matched into the path from $\text{parent}(\theta(q'))$ to $\text{root}(D)$ by Lemma 2, (iii) there are no right siblings on the path $\langle \theta(q') \cdots \text{prevSib}(d+1) \rangle$, because $\text{parent}(d+1)$ is an ancestor of $\theta(q')$, which is an ancestor of d , and (iv) $\text{subtree}(q')$ cannot be matched into $\text{subtree}(\hat{d})$ for every left sibling \hat{d} of the path from $\theta(q')$ to $\text{root}(D)$ by the induction hypothesis. From (i–iv) it follows that we cannot match $\text{subtree}(q')$ into the subtree rooted at $\text{prevSib}(d+1)$.

(Case 2): Otherwise, the next function call is $\text{L-MATCH}(d+1, q')$ in line 21. $\text{BACKTRACK}(d, q')$ has calculated the highest ancestor d' of the data node d such that $\langle d' \cdots \text{root}(D) \rangle$ matches $\langle \text{parent}(q') \cdots \text{root}(Q) \rangle$. By Lemma 2, d' equals $\theta(\text{parent}(q'))$.

We know that q' is the lowest node on $\langle q \cdots \text{root}(Q) \rangle$ such that $\theta(\text{parent}(q'))$ is an ancestor of $d+1$, because of the condition in the while loop. It follows that q' is matched somewhere on the path $\langle d \cdots \text{parent}(d+1) \rangle$ (for the case $q' \neq q$). No matter whether $q' = q$ or not, there was a function call $\text{L-MATCH}(q', d_0)$ for some d_0 on the path $\langle d \cdots \text{parent}(d+1) \rangle$. By the induction hypothesis and Lemma 2 there is no left sibling \hat{d} on the path $\langle d_0 \cdots \theta(\text{parent}(q')) \rangle$ such that $\text{subtree}(\hat{d})$ matches $\text{subtree}(q')$.

Since d_0 is in the subtree rooted at $\text{prevSib}(d+1)$, we now only need to show that $\text{subtree}(\text{prevSib}(d+1))$ does not match $\text{subtree}(q')$. We consider two cases:

- If $q = q'$, then we know that (i) $\text{subtree}(q)$ cannot be matched into $\text{subtree}(d)$, (ii) q cannot be matched into the path $\langle \text{parent}(d) \cdots \theta(\text{parent}(q)) \rangle$ by Lemma 2, (iii) there are no right siblings on the path $\langle d \cdots \text{prevSib}(d+1) \rangle$ due to the definition of the left-to-right pre-order successor, and (iv) $\text{subtree}(q)$ cannot be matched into $\text{subtree}(\hat{d})$ for every left sibling \hat{d} of the path $\langle d \cdots \theta(\text{parent}(q)) \rangle$

by the induction hypothesis. From (i-iv) it follows that the subtree rooted at $\text{prevSib}(d+1)$ does not match $\text{subtree}(q')$.

- If $q \neq q'$, there must have been a while loop execution with q' 's child on the path from q to q' as current query node and there must have been a function call $\text{L-MATCH}(\theta(q'), q')$. We know that (i) $\text{subtree}(q')$ cannot be matched into $\text{subtree}(\theta(q'))$ by Invariant 2, (ii) q' cannot be matched into the path $\langle \text{parent}(\theta(q')) \cdots \theta(\text{parent}(q)) \rangle$ by Lemma 2, (iii) there are no right siblings on the path $\langle \theta(q') \cdots \text{prevSib}(d+1) \rangle$, because there are no right siblings on the path $\langle d \cdots \text{prevSib}(d+1) \rangle$, and the path $\langle \theta(q') \cdots \text{prevSib}(d+1) \rangle$ is a subpath of that path (otherwise we would have returned a function call earlier, when the child of q' was the current query node), and (iv) $\text{subtree}(q')$ cannot be matched into $\text{subtree}(\hat{d})$ for every previous sibling \hat{d} of the path $\langle \theta(q') \cdots \theta(\text{parent}(q)) \rangle$ by the induction hypothesis. From (i-iv) it follows that we cannot match $\text{subtree}(q')$ into the subtree rooted at $\text{prevSib}(d+1)$.

- Otherwise there does not follow a function call. □

Proposition 4. *Algorithm 2 is complete. That is, given a data D and query tree Q , if Algorithm 2 returns false, then $D \not\equiv Q$.*

PROOF. We prove the proposition by induction on the number of nodes in the data tree D . If $|D| = 1$ then $\text{L-MATCH}(\text{root}(D), \text{root}(Q))$ returns true if $\text{root}(Q)$ is a leaf with an appropriate label in line 8 and false otherwise in line 15, which proves the completeness for that case.

Now suppose that $|D| > 1$. Assume L-MATCH returns false in line 15. Let d and q be the nodes such that, in the execution of $\text{L-MATCH}(d, q)$, false was returned. Due to line 13, d is a leaf and either q has children or the labels of q and d do not agree. Due to line 14, d is the maximal node w.r.t. $<_{\text{pre}}$, which means that there are no right siblings on the path from d to the root.

Consider a slight modification of the data tree: We attach an extra rightmost child to the root. Its value in the left-to-right pre-order is now $d+1$, the highest value of nodes in the data tree. Call this tree D' . Observe from the algorithm, that replacing D by D' does not make any difference in the function calls before $\text{L-MATCH}(d, q)$, because the algorithm traverses the data tree according to the left-to-right pre-order. However, in the function call $\text{L-MATCH}(d, q)$ the algorithm would not return false anymore, instead it would call $\text{L-MATCH}(d+1, q')$ for some query node q' . By Lemma 3 we know that for every child d' of the data root in D , $\text{subtree}(d')$ cannot match $\text{subtree}(q')$. We consider two cases.

- Assume that q' has a parent. It is clear that if there was a matching from Q into D , we would be able to match the $\text{subtree}(q')$ into some subtree of the data root. But we are not able to do this, so $D \not\equiv Q$.
- Assume that q' is the query root. By Lemma 2 we know that we cannot match the query root into the path $\langle \text{parent}(d+1) \cdots \text{root}(D) \rangle$. Hence, the labels of the query root and the data root do not agree and if there was a matching from Q into D , we would be able to match $\text{subtree}(q')$ into some subtree of the data root. But we are not able to do this, so $D \not\equiv Q$. □

Termination. Before we can conclude that L-MATCH is correct, we need to prove that the function call L-MATCH($\text{root}(D)$, $\text{root}(Q)$) terminates on every input D and Q . First, note that the while loop in line 18 terminates, because in every execution q' is overwritten with $\text{parent}(q')$ and our input trees are of finite depth.

We now only need to argue that whenever we call L-MATCH(d , q) for some $d \in D$ and $q \in Q$, we have not called L-MATCH(d , q) before. We prove this in the following lemma (it is an immediate consequence of Lemma 4 letting $q_0 = q$ and $d_0 = d$).

Lemma 4. *Let L-MATCH(d , q) be a function call resulting from the initial procedure call L-MATCH($\text{root}(D)$, $\text{root}(Q)$). Then at the time when L-MATCH(d , q) is called*

$$\forall q_0 \geq q, \forall d_0 \geq d, \text{ we have not yet called L-MATCH}(d_0, q_0) \text{ before.}$$

PROOF. We prove the lemma by induction on the position k of L-MATCH(d , q) in the sequence of function calls resulting from the initial procedure call.

More specifically, our induction hypothesis will be: at the time when L-MATCH(d , q) is called,

(I1): $\forall q_0 \geq q, \forall d_0 \geq d$, we have not yet called L-MATCH(d_0 , q_0) before;

(I2): for all right siblings \hat{q} of nodes on the path $\langle q \cdots \text{root}(Q) \rangle$, for all nodes $q_0 \in \text{subtree}(\hat{q})$, and for all data nodes $d_0 \in \text{subtree}(\theta(\text{parent}(\hat{q})))$, we have not yet called L-MATCH(d_0 , q_0).

(I3): for all nodes $\hat{q} < q$, for all nodes $q_0 \in \text{subtree}(\hat{q})$, for all right siblings \hat{d} on the path $\langle \theta(\hat{q}) \cdots \text{root}(D) \rangle$, and for all data nodes $d_0 \in \text{subtree}(\hat{d})$, we have not yet called L-MATCH(d_0 , q_0).

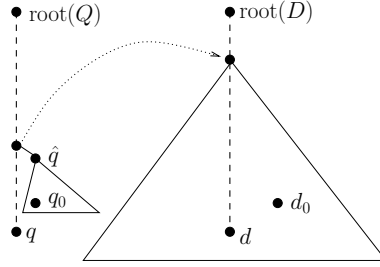
We illustrate the hypotheses (I2) and (I3) in Figure 2.

If $k = 1$ then we have L-MATCH($\text{root}(D)$, $\text{root}(Q)$) in which case there is nothing to show.

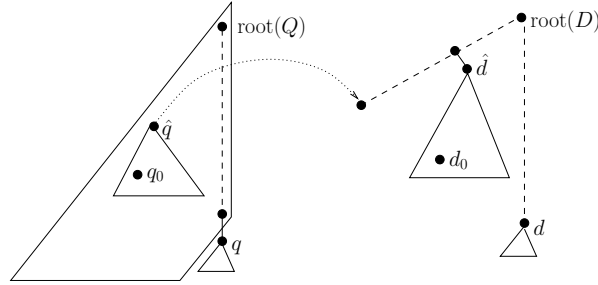
So, from now on we assume that the Lemma holds for the first k function calls. Let L-MATCH(d , q) be the k^{th} function call. We prove that the Lemma also holds for the $k + 1^{\text{th}}$ function call (if there is one).

Let us start with a simple observation concerning (I3). The induction hypothesis for (I3) implies that, for all query nodes $\hat{q} < q$, for all query nodes $q_0 \in \text{subtree}(\hat{q})$, for all right siblings \hat{d} on the path $\langle \theta(\hat{q}) \cdots \text{root}(D) \rangle$, and for all data nodes $d_0 \in \text{subtree}(\hat{d})$ we have not called L-MATCH(d_0 , q_0) before we called L-MATCH(d , q). We argue why this remains true even after calling L-MATCH(d , q), but before the next function call is made. Towards a contradiction, assume that this was not the case. In that case there would be a query node $\hat{q} < q$ such that $q \in \text{subtree}(\hat{q})$, and a right sibling \hat{d} on the path $\langle \theta(\hat{q}) \cdots \text{root}(D) \rangle$, such that $d \in \text{subtree}(\hat{d})$. But this cannot be because, due to Lemma 2, the ancestors of q are matched on the path $\langle d \cdots \text{root}(D) \rangle$ and hence d cannot be in a subtree of a right sibling on the path $\langle \theta(\hat{q}) \cdots \text{root}(D) \rangle$. Hence, (I3) is also still true right after calling L-MATCH(d , q). (\dagger)

Next we will consider the four possible function calls following the k^{th} function call L-MATCH(d , q) and we will show that (I1), (I2), and (I3) still hold for the next function call.



(a) Induction hypothesis (I2).



(b) Induction hypothesis (I3).

Figure 2: Illustrations of the induction hypotheses in the proof of Lemma 4.

- If the next function call is $\text{L-MATCH}(d+1, q+1)$ (line 3), then $\theta(q)$ is defined to be d . Here, $d+1$ is the leftmost child of d and $q+1$ is the leftmost child of q . The induction hypothesis for (I1) implies that for all $q_0 \geq q$, for all data nodes $d_0 \geq d$, we have not called $\text{L-MATCH}(d_0, q_0)$ before we called $\text{L-MATCH}(d, q)$. In the meantime, we only executed $\text{L-MATCH}(d, q)$, so for all $q_0 \geq q+1$, for all data nodes $d_0 \geq d+1$, we have not called $\text{L-MATCH}(d_0, q_0)$, which proves (I1).

Item (I2) of the induction hypothesis implies that, for all right siblings \hat{q} of nodes on the path $\langle q \cdots \text{root}(Q) \rangle$, for all nodes $q_0 \in \text{subtree}(\hat{q})$, for all data nodes $d_0 \in \text{subtree}(\theta(\text{parent}(\hat{q})))$, we have not called $\text{L-MATCH}(d_0, q_0)$ before we called $\text{L-MATCH}(d, q)$. Since $q+1$ is the leftmost child of q , we only need to show that, for all right siblings \hat{q} of $q+1$, for all nodes $q_0 \in \text{subtree}(\hat{q})$, for all data nodes $d_0 \in \text{subtree}(\theta(\text{parent}(\hat{q})))$ (which is the $\text{subtree}(d)$), we have not called $\text{L-MATCH}(d_0, q_0)$ before. But this follows from the induction on (I1), because data nodes in $\text{subtree}(d)$ are greater or equal to d and query nodes in subtrees of $q+1$'s right siblings are greater than q . This shows (I2).

In order to prove (I3) we need to show that for $\hat{q} < q+1$, for all query nodes $q_0 \in \text{subtree}(\hat{q})$, for all right siblings \hat{d} on the path $\langle \theta(\hat{q}) \cdots \text{root}(D) \rangle$, and for all data nodes $d_0 \in \text{subtree}(\hat{d})$ we have not called $\text{L-MATCH}(d_0, q_0)$ before calling $\text{L-MATCH}(d+1, q+1)$.

By the observation (†) above this is true for $\hat{q} < q$. So, let us consider $\hat{q} = q$. The fact that $\theta(q) = d$ now implies that $d_0 > d$ and $q_0 \geq q$. The claim follows from the induction on item (1) and the fact that we only called $\text{L-MATCH}(d, q)$ in the

meantime.

- If the next function call is $\text{L-MATCH}(d+1, q)$ (line 5), then $d+1$ is the leftmost child of d . The induction on (I1) implies that (I1) is true (as above). Since we only called $\text{L-MATCH}(d, q)$ in the meantime and did not change the mapping θ at all, (I2) is a direct consequence of the induction hypothesis on (I2). Since we did not change the mapping θ and the query node serving as argument of the $k+1^{\text{th}}$ function call is the same as the argument of the k^{th} function call, (I3) follows from the observation (\dagger) made above.
- If the next function call is $\text{L-MATCH}(d'+1, q+1)$ (line 11), then $\theta(q)$ is defined to be d . Here, $q+1$ is a right sibling of a node on the path $\langle q \cdots \text{root}(Q) \rangle$ (due to the left-to-right pre-order and the fact that q is a leaf, see line 6) and $d'+1$ is the leftmost child of some ancestor of d . The induction on (I1) assures that, for all $q_0 \geq q$, for all data nodes $d_0 \geq d$, we have not called $\text{L-MATCH}(d_0, q_0)$ before we called $\text{L-MATCH}(d, q)$. In the meantime, we executed $\text{L-MATCH}(d, q)$, so for all $q_0 \geq q+1$, for all data nodes $d_0 \geq d$, we have not called $\text{L-MATCH}(d_0, q_0)$. In order to prove (I1), we still need to show that this is also true for all $q_0 \geq q+1$ and for all data nodes d_0 with $d'+1 \leq d_0 < d$. We consider two cases:
 - If $q_0 \in \text{subtree}(q+1)$ we can make use of the induction hypothesis on (I2). The query node $q+1$ is a right sibling of a node on the path $\langle q \cdots \text{root}(Q) \rangle$ and hence we have not called $\text{L-MATCH}(d_0, q_0)$ before for all nodes $d_0 \in \text{subtree}(\theta(\text{parent}(q+1)))$. This proves our case because, by Lemma 2, $\theta(\text{parent}(q+1))$ is equal to d' and d' is an ancestor of d . Clearly, $\text{subtree}(d')$ includes all nodes d_0 with $d'+1 \leq d_0 < d$. Hence, for all $q_0 \in \text{subtree}(q+1)$ and for all d_0 with $d'+1 \leq d_0 < d$ we have not called $\text{L-MATCH}(d_0, q_0)$ before.
 - If $q_0 \notin \text{subtree}(q+1)$ we can make use of the induction hypothesis on (I2) again. By definition of the left-to-right pre-order, q_0 is then in a subtree of some right sibling \hat{q} of a node on the path $\langle q+1 \cdots \text{root}(Q) \rangle$. This \hat{q} is also a right sibling of a node on the path $\langle q \cdots \text{root}(Q) \rangle$. By induction on (I2) it follows that, for all data nodes $d_0 \in \text{subtree}(\theta(\text{parent}(\hat{q})))$, we have not called $\text{L-MATCH}(d_0, q_0)$. This proves our case, because $\text{parent}(\hat{q})$ is an ancestor of or equal to $\text{parent}(q+1)$ and, by Lemma 2, $\theta(\text{parent}(\hat{q}))$ is an ancestor of or equal to $\theta(\text{parent}(q+1))$, which is equal to d' . Clearly, $\text{subtree}(d')$ includes all nodes d_0 with $d'+1 \leq d_0 < d$ and so does $\text{subtree}(\theta(\text{parent}(\hat{q})))$. Hence, for all $q_0 \notin \text{subtree}(q+1)$ and for all d_0 with $d'+1 \leq d_0 < d$ we have not called $\text{L-MATCH}(d_0, q_0)$ before.

As mentioned above, right siblings of a node on the path $\langle q+1 \cdots \text{root}(Q) \rangle$ are also a right siblings of a node on the path $\langle q \cdots \text{root}(Q) \rangle$. Hence, (I2) immediately follows from the induction hypothesis on (I2).

In order to prove (I3) we need to show that, for $\hat{q} < q+1$, for all query nodes $q_0 \in \text{subtree}(\hat{q})$, for all right siblings \hat{d} on the path $\langle \theta(\hat{q}) \cdots \text{root}(D) \rangle$, and for all data nodes $d_0 \in \text{subtree}(\hat{d})$, we have not called $\text{L-MATCH}(d_0, q_0)$ before calling $\text{L-MATCH}(d'+1, q+1)$.

By the observation (\dagger) above this is true for $\hat{q} < q$. So, let us consider $\hat{q} = q$. The left-to-right pre-order and the fact that $\theta(q) = d$, implies that $d_0 > d$ and

$q_0 \geq q$. The claim follows from the induction on (I1) and the fact that we only called $\text{L-MATCH}(d, q)$ in the meantime.

- If the next function call is $\text{L-MATCH}(d + 1, q')$ (line 21 or line 25), then $d + 1$ is a right sibling of a node on the path $\langle d \cdots \text{root}(D) \rangle$ (due to the left-to-right pre-order and the fact that d is a leaf, see line 13) and q' is an ancestor of or is equal to q .

The induction on (I1) implies that, for all $q_0 \geq q$, for all data nodes $d_0 \geq d$, we have not called $\text{L-MATCH}(d_0, q_0)$ before we called $\text{L-MATCH}(d, q)$. In the meantime, we only executed $\text{L-MATCH}(d, q)$, so, for all $q_0 \geq q$, for all data nodes $d_0 \geq d + 1$, we have not called $\text{L-MATCH}(d_0, q_0)$ before.

In order to prove (I1) we still need to show that this is also true for all query nodes q_0 with $q' \leq q_0 < q$ and for all data nodes $d_0 \geq d + 1$. So, take a query node q_0 such that $q' \leq q_0 < q$. Note that such a query node q_0 is in subtree(q'). Furthermore, each node d_0 that is greater or equal to $d + 1$ is in the subtree of some right sibling \hat{d} on the path $\langle \text{prevSib}(d + 1) \cdots \text{root}(D) \rangle$ because d is a leaf. This path is a subpath of $\langle \theta(q') \cdots \text{root}(D) \rangle$, because q' is the lowest ancestor of q whose parent is an ancestor of $d + 1$, which means that q' is mapped onto a node on $\langle \text{parent}(d) \cdots \text{prevSib}(d + 1) \rangle$ by Lemma 2 and the fact that $q' < q$. By induction on (I3) it follows that we have not called $\text{L-MATCH}(d_0, q_0)$ before calling $\text{L-MATCH}(d, q)$.

Since q' is an ancestor of or equal to q , the path $\langle q' \cdots \text{root}(Q) \rangle$ is a sub-path of the path $\langle q \cdots \text{root}(Q) \rangle$. Hence, (I2) immediately follows from the induction on (I2).

Since we did not change the mapping θ and the query node serving as argument of the $k + 1^{\text{th}}$ function call is smaller or equal to the argument of the k^{th} function call, (I3) follows from the observation (\dagger) made above. \square

Propositions 3, 4, and Lemma 4 imply the correctness of L-MATCH .

Proposition 5. *Algorithm 2 is correct. That is, given the roots d and q of a data D and query tree Q , $\text{L-MATCH}(d, q)$ decides whether $D \models Q$.*

3.2.2. Space Complexity of L-MATCH .

We already argued in the main body of the paper that the recursion stack has no influence on the operation of L-MATCH . It remains to argue why BACKTRACK only needs logarithmic space. $\text{BACKTRACK}(d, q)$ calculates the highest ancestor d' of the data node d such that the path $\langle d' \cdots \text{root}(D) \rangle$ matches the path $\langle \text{parent}(q) \cdots \text{root}(Q) \rangle$. The difficulty lies in the fact that we cannot store both paths. Instead, we store d and q . We also store two help variables d_0 and q_0 , which are initialized to be $\text{root}(D)$ and $\text{root}(Q)$, respectively. We now iterate over the following. We compare the labels of d_0 and q_0 . If they match, we overwrite d_0 and q_0 with the children of d_0 and q_0 that lie on the paths to d and q , respectively. This is performed as explained in the beginning of this section: We can start at d (resp., q), scan the input tape for the unique node that has a child pointer to d (resp., q), and continue upwards in this manner until we find a child of d_0 , resp., q_0 . If the labels of d_0 and q_0 do not match, we only overwrite d_0 with its child on the path to d . We continue until we matched the whole path $\langle \text{parent}(q) \cdots \text{root}(Q) \rangle$. Finally we return the data node onto which we matched $\text{parent}(q)$.

3.2.3. The Complexity of the Tree Homeomorphism Problem.

As argued above, L-MATCH can be performed in LOGSPACE. Putting this together with the fact that reachability in trees is LOGSPACE-complete, given the tree as a pointer structure [8], we obtain the following Theorem.

Theorem 3. *The tree homeomorphism problem is LOGSPACE-complete.*

4. The Bottom-up Algorithm

Although the previously presented top-down algorithms for tree homeomorphism matching are quite space-efficient, their time complexity is quite high and they involve quite a lot of recomputing of already obtained matchings, which is unsatisfactory. We therefore turn to a bottom-up matching approach which has the property that *no* obtained matchings between the data and query tree need to be recomputed, which leads to a better time complexity of the overall algorithm.

Before presenting the bottom-up algorithm for the tree homeomorphism matching problem in detail, we need to introduce several formal notions. As in the previous section, we first present an algorithm for the tree homeomorphism problem and then show how to change it into an algorithm for the tree homeomorphism matching problem.

In the present section, we assume the *left-to-right post-order* ordering $<_{\text{post}}$ on nodes in trees and hedges. For a node u , we denote by $u + 1$ and $u - 1$ the successor and predecessor of u in the left-to-right postorder ordering, respectively. Moreover, when we, e.g., use terminology such as “largest” and “smallest”, we always assume the left-to-right post ordering. In this section, we also assume that XML documents are stored on tape in left-to-right postorder (or, alternatively, together with a left-to-right postorder index), which allows a random-access machine model to verify the left-to-right post-order ordering in constant time. To simplify the presentation of our algorithm, we also assume two dummy nodes in every tree and hedge: nil and ∞ . The node nil is such that $\text{nil} + 1$ is the smallest node in the hedge, and the node ∞ is defined as the successor of the largest node of the hedge. Given two nodes $h_{\text{from}} \leq h_{\text{until}}$ in a hedge H , we denote by the interval $[h_{\text{from}}, h_{\text{until}}]$ the subhedge of H consisting only of the nodes $\{v \mid h_{\text{from}} \leq v \leq h_{\text{until}}\}$.⁴ The notion of such an interval in a tree is illustrated in Figure 3(a). Here, the interval $[h_{\text{from}}, h_{\text{until}}]$ is the striped area in the tree. Given a hedge H and a node $h \in \text{Nodes}(H)$, we denote by $\text{subhedge}^H(h)$ the subhedge $[h_{\text{from}}, h]$, where h_{from} is the smallest descendant of h 's leftmost sibling according to the left-to-right postorder ordering. We illustrate this notion in Figure 3(b).

When H is a data hedge or a tree pattern query, we refer to $[h_{\text{from}}, h_{\text{until}}]$ as a data or query hedge interval, respectively. We extend the semantics of tree pattern matching to hedges as follows. Let $Q_1 \cdots Q_n$ be a query hedge interval $[q_{\text{from}}, q_{\text{until}}]$ and $D_1 \cdots D_m$ be a data hedge interval $[d_{\text{from}}, d_{\text{until}}]$. We say that $[d_{\text{from}}, d_{\text{until}}]$ matches $[q_{\text{from}}, q_{\text{until}}]$, denoted by $[d_{\text{from}}, d_{\text{until}}] \models [q_{\text{from}}, q_{\text{until}}]$, if, for every Q_i , $i = 1, \dots, n$, there exists a D_j , $j = 1, \dots, m$, such that $D_j \models Q_i$.

Before presenting the intuition about the bottom-up tree homeomorphism algorithm, we describe an auxiliary procedure RTOP, which, given two nodes h_{from} and h_{until} ,

⁴Notice that our definition of a hedge did not assume all root nodes of the individual trees to be siblings of one another.

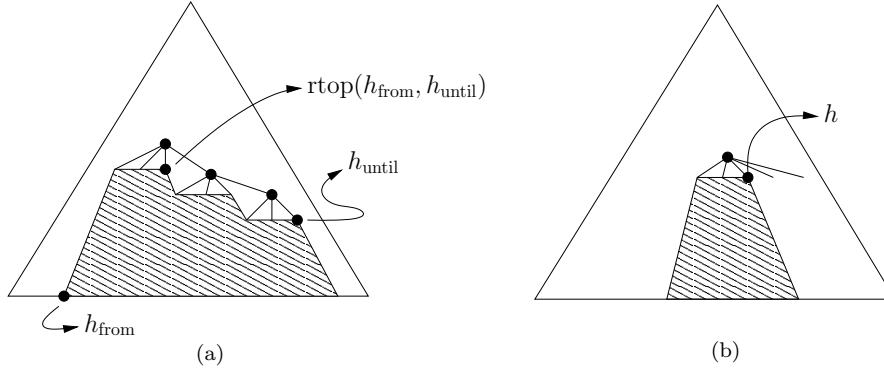


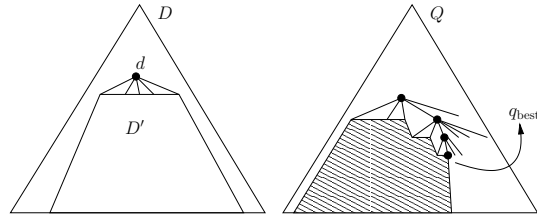
Figure 3: Illustration of a hedge interval and RTOP (left) and of subhedge^H(h) (right).

returns the rightmost node among the topmost nodes in the interval $[h_{\text{from}}, h_{\text{until}}]$. More formally, $\text{RTOP}(h_{\text{from}}, h_{\text{until}})$ is the node u such that $\text{depth}(u)$ is minimal and u is larger than every other node v in $[h_{\text{from}}, h_{\text{until}}]$ with $\text{depth}(u) = \text{depth}(v)$. This notion is illustrated in Figure 3(a). Furthermore, in order to simplify the presentation of the algorithm, we define $\text{RTOP}(h_{\text{from}}, h_{\text{until}}) = \infty$ if $h_{\text{from}} > h_{\text{until}}$. Notice that RTOP can easily be computed in time linear in the depth of the tree and in logarithmic space by traversing the path from h_{until} to the query root and comparing the previous siblings of nodes on the path with h_{from} w.r.t. the left-to-right post-ordering. Indeed, assume that $h_{\text{from}} \leq h_{\text{until}}$. Let u be the highest ancestor of h_{until} that has a previous sibling s such that $s \geq h_{\text{from}}$. If no such u exists, then $\text{rtop}(h_{\text{from}}, h_{\text{until}})$ is h_{until} . Otherwise, $\text{rtop}(h_{\text{from}}, h_{\text{until}})$ is s .

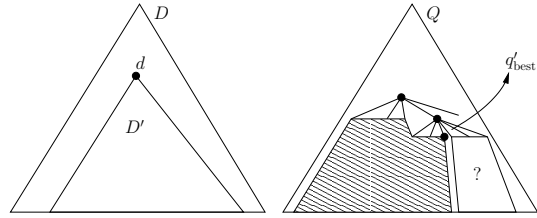
We first present an algorithm for *deciding* whether $D \models Q$ and show later how it can be extended to an algorithm for the tree homeomorphism matching problem. The main procedure of our algorithm is called TMATCH. Given a data node d and query nodes q_{from} and q_{until} , TMATCH returns the largest query node q in the interval $[q_{\text{from}}, q_{\text{until}}]$ such that $\text{subtree}^D(d)$ matches $[q_{\text{from}}, q]$ if q exists; and $q_{\text{from}} - 1$ otherwise. Hence, if d is the root of D , and q_{from} and q_{until} are the leftmost leaf and the root of Q , respectively, then $D \models Q$ if and only if TMATCH returns q_{until} .

TMATCH uses an auxiliary procedure called HMATCH, which, given a data node d and query nodes q_{from} and q_{until} , returns the largest node q in the interval $[q_{\text{from}}, q_{\text{until}}]$ such that $\text{subhedge}^D(d)$ matches $[q_{\text{from}}, q]$ if q exists; and $q_{\text{from}} - 1$ otherwise.

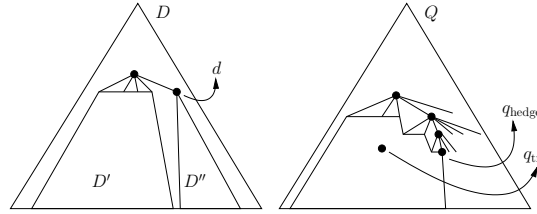
We start by explaining the operation of TMATCH, which is presented in Algorithm 3. Given a data node d and query nodes q_{from} and q_{until} , TMATCH first starts by recursively calling HMATCH with the same query nodes for the subhedge D' of D defined by d 's last child, yielding result q_{best} (see Figure 4(a)). In the remainder of TMATCH, we essentially want to test how q_{best} can be improved when we also consider the node d in addition to D' . One particular interesting case is when q_{best} is a last sibling and its parent has the same label as d . In this case, we can at least improve our best query node to q_{best} 's parent which we call here q'_{best} . Furthermore, it is possible that q'_{best} is not yet the best query node we can obtain. In particular, we still need to test which part of the hedge defined by $[q'_{\text{best}} + 1, \text{lastSib}(q'_{\text{best}})]$ can be matched in the subtree below d (see Figure 4(b)). The



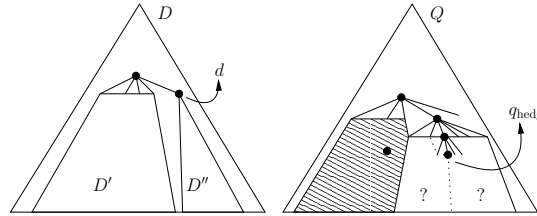
(a) Operation of TMATCH: recursive call of HMATCH.



(b) Operation of TMATCH: recursive call of TMATCH.



(c) Operation of HMATCH: first recursive calls of TMATCH and HMATCH.



(d) Operation of HMATCH: a subsequent recursive call of TMATCH, trying to improve q_{tree} .

Figure 4: Illustrations of the tree homeomorphism algorithm.

largest node that is obtained in this manner is the node that TMATCH should return.

We now explain the operation of HMATCH, which is presented in Algorithm 4. Essentially, given d , q_{from} , and q_{until} , HMATCH starts by recursively calling itself with the same query nodes on the hedge defined by the previous sibling of d (i.e., D' in Figure 4(c)), yielding q_{hedge} , and by calling TMATCH with the same query nodes on the subtree under d itself (D'' in Figure 4(c)), yielding q_{tree} . The remainder of HMATCH consists of iteratively improving q_{tree} and q_{hedge} . That is, while it is possible that D' and D'' yield

Algorithm 3 Function TMATCH. Here, $+1$ and -1 denote the successor and predecessor in the depth-first left-to-right post-ordering, respectively.

```

TMATCH (DNode  $d$ , QNode  $q_{\text{from}}$ , QNode  $q_{\text{until}}$ )
2: if  $d$  is a leaf then  $q_{\text{best}} \leftarrow q_{\text{from}} - 1$ 
   else  $q_{\text{best}} \leftarrow \text{HMATCH}(\text{lastChild}(d), q_{\text{from}}, q_{\text{until}})$ 
4: end if
   if  $q_{\text{best}} + 1 \leq_{\text{post}} q_{\text{until}}$  and  $d$  matches  $q_{\text{best}} + 1$  then
6:    $q_{\text{best}} \leftarrow q_{\text{best}} + 1$ 
     if  $q_{\text{best}} + 1 \leq_{\text{post}} \text{lastSib}(q_{\text{best}})$  then
8:   return TMATCH( $d, q_{\text{best}} + 1, \text{lastSib}(q_{\text{best}})$ )
     else return  $q_{\text{best}}$ 
10:  end if
   else return  $q_{\text{best}}$ 
12: end if

```

small values of q_{tree} and q_{hedge} , their concatenation can give rise to a much larger part of the query that can be matched. Essentially, this is due to the fact that the matching of tree pattern queries is *unordered*. For example, it can occur that we need to match a certain first sibling in D' , a second one in D'' , a third one again in D' and so on. Hence, the procedure HMATCH alternates between finding best matches in D' and D'' until it reaches a fixpoint.

However, we need to take care in how this fixpoint is computed. One possible case is illustrated in Figure 4(d). This particular case builds further on the situation in Figure 4(c). Here, we try to improve q_{tree} by starting the TMATCH procedure again for the node d , but now only with the part of the query marked with question marks. The case where q_{tree} is larger than q_{hedge} is dual and not illustrated here.

Example 1. Figure 5(a) and 5(b) illustrate an example for the bottom up algorithm. For brevity, we denote TMATCH and HMATCH with TM and HM, respectively. The first calls of TM and HM demonstrate the basic recursive structure of our algorithm: TM on a node d calls HM on the rightmost child of d . HM on a node d returns TM of d if that node is a first sibling; or performs a divide-and-conquer technique by calling HM on the left sibling of d and TM on d itself (as in the function call $\text{HM}(d_4, q_1, q_5)$). Further recursive calls to TM or HM are then needed to maximize the part of the query that can be matched.

The simplest function call in the example that performs such further recursive calls is the call $\text{HM}(d_2, q_1, q_5)$, which starts by computing $q_{\text{hedge}} = \text{HM}(d_1, q_1, q_5)$ and $q_{\text{tree}} = \text{TM}(d_2, q_1, q_5)$. As can be seen in Figure 5(b), $q_{\text{hedge}} = \text{nil}$. The call $\text{TM}(d_2, q_1, q_5)$ is more successful, because d_2 and q_1 are both labeled with a . In general, it might be possible that q_2 and further nodes can be matched in $\text{subtree}(d_2)$. The function call $\text{TM}(d_2, q_2, q_4)$ checks that possibility. (For sure, q_1 and q_5 cannot both be matched on d_2 , which is why we restrict the query tree interval by q_4 .) But q_2 is not labeled with a so the return value of the two TM calls is q_1 . After this initial phase, $\text{HM}(d_2, q_1, q_5)$ tries to improve q_{tree} and q_{hedge} iteratively. It calls $\text{HM}(d_1, q_2, q_4)$ and improves q_{hedge} to be q_2 , because q_2 and d_1 are both labeled with b . Further improvements fail as there is no c -labeled node in the subhedge of d_2 .

Algorithm 4 Function HMATCH. Here, $+1$ and -1 denote the successor and predecessor in the depth-first left-to-right post-ordering, respectively.

```

HMATCH (DNode  $d$ , QNode  $q_{\text{from}}$ , QNode  $q_{\text{until}}$ )
2: if  $d$  is a first sibling then return TMATCH( $d$ ,  $q_{\text{from}}$ ,  $q_{\text{until}}$ )
   else
4:    $q_{\text{hedge}} \leftarrow$  HMATCH(prevSib( $d$ ),  $q_{\text{from}}$ ,  $q_{\text{until}}$ )
      $q_{\text{tree}} \leftarrow$  TMATCH( $d$ ,  $q_{\text{from}}$ ,  $q_{\text{until}}$ )
6:   loop
     if  $q_{\text{hedge}} = q_{\text{tree}}$  then return  $q_{\text{hedge}}$ 
8:   else if  $q_{\text{tree}} <_{\text{post}} q_{\text{hedge}}$  then
      $r_{\text{top}} \leftarrow$  RTOP( $q_{\text{tree}} + 1$ ,  $q_{\text{hedge}}$ )
10:    while  $r_{\text{top}} <_{\text{post}} \infty$  and  $q_{\text{hedge}} <_{\text{post}} \text{lastSib}(r_{\text{top}})$  do
      $q_{\text{tree}} \leftarrow$  TMATCH( $d$ ,  $r_{\text{top}}+1$ ,  $\text{lastSib}(r_{\text{top}})$ )
12:     $r_{\text{top}} \leftarrow$  RTOP( $q_{\text{tree}} + 1$ ,  $q_{\text{hedge}}$ )
     end while
14:    if  $q_{\text{tree}} \leq_{\text{post}} q_{\text{hedge}}$  then return  $q_{\text{hedge}}$ 
     end if
16:   else
      $r_{\text{top}} \leftarrow$  RTOP( $q_{\text{hedge}} + 1$ ,  $q_{\text{tree}}$ )
18:    while  $r_{\text{top}} <_{\text{post}} \infty$  and  $q_{\text{tree}} <_{\text{post}} \text{lastSib}(r_{\text{top}})$  do
      $q_{\text{hedge}} \leftarrow$  HMATCH(prevSib( $d$ ),  $r_{\text{top}} + 1$ ,  $\text{lastSib}(r_{\text{top}})$ )
20:     $r_{\text{top}} \leftarrow$  RTOP( $q_{\text{hedge}} + 1$ ,  $q_{\text{tree}}$ )
     end while
22:    if  $q_{\text{hedge}} \leq_{\text{post}} q_{\text{tree}}$  then return  $q_{\text{tree}}$ 
     end if
24:   end if
   end loop
26: end if

```

A similar iterative improvement is illustrated by $\text{HM}(d_3, q_1, q_5)$. Observe that we try to improve q_{tree} here and call $\text{TM}(d_4, q_2, q_4)$ and $\text{TM}(d_4, q_3, q_3)$. Only the latter call yields an improvement. But we cannot omit the former one: if $\text{subtree}(d_4)$ would match $\text{subtree}(q_4)$, then the former call would yield q_4 and the latter call would yield q_3 . As we want our algorithm to return the largest query node such that the interval ending with it can be matched the result of the former call would have been the relevant one in that case.

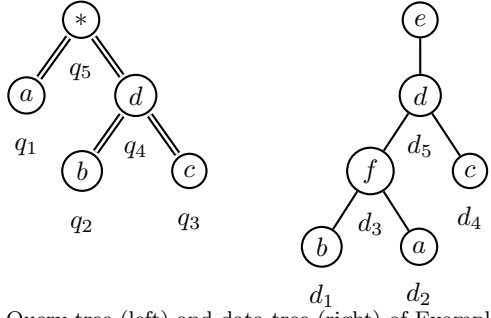
4.1. Correctness.

The main technical difficulty of this section is proving that TMATCH is correct.

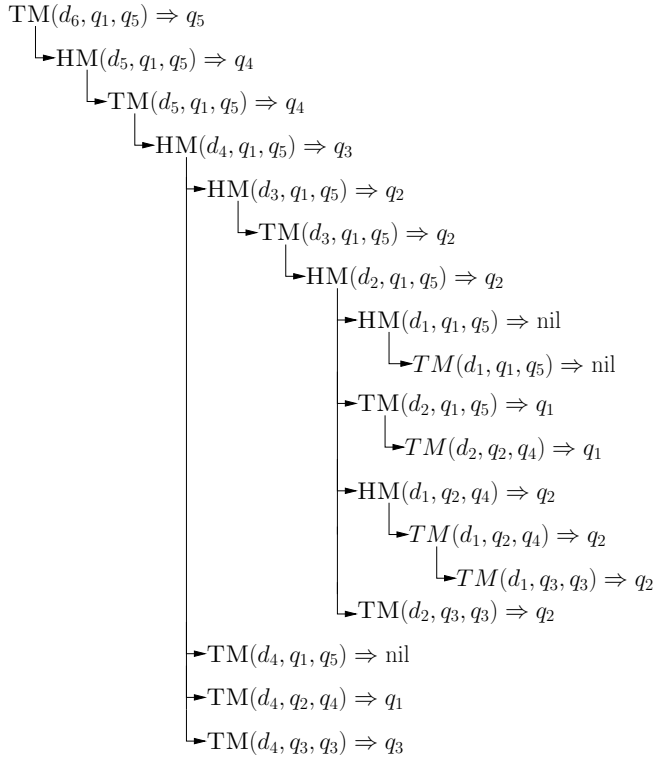
Lemma 5. *Let D be a data tree and let Q be a query tree. TMATCH is correct, that is, given the root node d of D , the smallest and largest node q_{from} and q_{until} of Q , respectively, TMATCH returns q_{until} iff $D \models Q$.*

For the proof of Lemma 5, we start with a few simple observations.

Observation 4. *A node u is not a last sibling $\Leftrightarrow u + 1$ is a leaf.*



(a) Query tree (left) and data tree (right) of Example 1.



(b) Function calls of HMATCH (HM) and TMATCH (TM) of Example 1.

Figure 5: Illustrations for Example 1.

PROOF. Left to right: if u is not a last sibling, then $u + 1$ is the leftmost descendant leaf of the right sibling of u , or the right sibling of u itself if it is a leaf. Right to left: if u is the last node in a left-to-right postorder traversal, then u is a last sibling for which $u + 1$ does not exist. For all other last siblings u , $u + 1$ is u 's parent, which is not a leaf. \square

We call a hedge interval *complete* when if it contains a certain node, it also contains its children.

Observation 5. *In Algorithms 3 and 4, the following properties hold:*

- (1) q_{until} is always a last sibling.
- (2) q_{from} is always a leaf.
- (3) $[q_{\text{from}}, q_{\text{until}}]$ is always a complete interval.

PROOF. (1) In our initial call of TMATCH, q_{until} is the root node of the tree, which is always a last sibling. The property for the deeper recursive calls follows immediately from a straightforward inspection of the recursive function calls in the algorithm.

(2) In our initial call of TMATCH, q_{from} is the smallest node of Q , which is always a leaf. Furthermore, in TMATCH we only call HMATCH with q_{from} as a second parameter and TMATCH with $q_{\text{best}} + 1$ as a second parameter if q_{best} is not a last sibling (which is a leaf due to Observation 4). In HMATCH all recursive calls have either q_{from} or $\text{rtop} + 1$ as second parameter. We show that, in this case, rtop is never a last sibling. Hence, according to Observation 4, $\text{rtop} + 1$ is always a leaf. In the calls of TMATCH on line 11, we have that $\text{rtop} < \infty$ and $q_{\text{hedge}} < \text{lastSib}(\text{rtop})$, due to the while condition. As $\text{rtop} < \infty$, we have that $\text{rtop} \leq q_{\text{hedge}}$ due to the calls of RTOP on line 9 and line 12. Hence, $\text{rtop} < \text{lastSib}(\text{rtop})$. The proof is analogous for the calls of HMATCH on line 19.

(3) In the initial call of TMATCH, the claim obviously holds. In TMATCH we call HMATCH with q_{from} and q_{until} , for which the claim then trivially also holds; and TMATCH with $q_{\text{best}} + 1$ and $\text{lastSib}(q_{\text{best}})$ if q_{best} is not a last sibling. Hence, $[q_{\text{best}} + 1, \text{lastSib}(q_{\text{best}})]$ is equal to the hedge subtree($\text{nextSib}(q_{\text{best}})$) \cdots subtree($\text{lastSib}(q_{\text{best}})$), which is complete. The proof for the recursive calls in HMATCH is analogous. \square

Observation 6. *Let d_1 and d_2 be data nodes and q be a query node. If $[d_1, d_2]$ does not match subtree(q), then $[d_1, d_2]$ does not match any query tree interval containing subtree(q).*

PROOF. Let q_{from} and q_{until} be such that $[q_{\text{from}}, q_{\text{until}}] = \text{subtree}(q)$. For $q'_{\text{from}} \leq q_{\text{from}}$ and $q'_{\text{until}} \geq q_{\text{until}}$, it can be shown by a simple structural induction on the hedge $[q'_{\text{from}}, q'_{\text{until}}]$ that $[d_1, d_2]$ does not match $[q'_{\text{from}}, q'_{\text{until}}]$. \square

Observation 7. *Let H be a data hedge and $[q_{\text{from}}, q_{\text{until}}]$ be a complete query tree interval. We have that q is the largest node in $[q_{\text{from}}, q_{\text{until}}]$ such that $H \models [q_{\text{from}}, q]$ if and only if*

- H matches $[q_{\text{from}}, q]$; and
- either $q = q_{\text{until}}$ or H does not match subtree($q + 1$).

PROOF. Left to right: Let H be a data hedge and let $[q_{\text{from}}, q_{\text{until}}]$ be a query tree interval. Let q be the largest node in $[q_{\text{from}}, q_{\text{until}}]$ such that $H \models [q_{\text{from}}, q]$. If $q = q_{\text{until}}$ we are done. Otherwise, if, towards a contradiction, H matches subtree($q + 1$), then we also immediately have that H matches $[q_{\text{from}}, q + 1]$, which contradicts the maximality of q .

Right to left: Let q be a query node in $[q_{\text{from}}, q_{\text{until}}]$ such that H matches $[q_{\text{from}}, q]$. If $q = q_{\text{until}}$ then we are done. Otherwise, notice that, as $q + 1$ is in the complete interval $[q_{\text{from}}, q_{\text{until}}]$, we have that subtree($q + 1$) is entirely contained in $[q_{\text{from}}, q_{\text{until}}]$. Hence, if H does not match subtree($q + 1$), then H also cannot match $[q_{\text{from}}, q + 1]$. The latter can be shown by a simple structural induction on $[q_{\text{from}}, q + 1]$. \square

4.1.1. Correctness of TMatch.

For readability, we split the correctness proof into several lemmas. Essentially, the proof is by induction on the height of the data node d in D .

Lemma 6. *Let d be a leaf data node and q_{from} and q_{until} be query nodes. Given d , q_{from} , and q_{until} , TMATCH is correct, that is, TMATCH returns the largest node q in $[q_{\text{from}}, q_{\text{until}}]$ such that $\text{subtree}(d) \models [q_{\text{from}}, q]$ if it exists; and $q_{\text{from}} - 1$ otherwise.*

PROOF. By induction on the number of nodes of $[q_{\text{from}}, q_{\text{until}}]$.

$q_{\text{from}} = q_{\text{until}}$: We initialize q_{best} with $q_{\text{from}} - 1$ on line 2. If d does not match q_{from} on line 5, we immediately return $q_{\text{best}} = q_{\text{from}} - 1$ on line 11. If d matches $q_{\text{from}} = q_{\text{until}}$ on line 5, q_{best} gets the value q_{from} on line 6. As $q_{\text{from}} = q_{\text{until}}$ is a last sibling (Observation 5), we do not execute the recursive call on line 8 and return q_{from} in line 9. Both cases are easily seen to be correct.

$q_{\text{from}} < q_{\text{until}}$: We initialize q_{best} with $q_{\text{from}} - 1$ on line 2. If d does not match q_{from} on line 5, we return $q_{\text{best}} = q_{\text{from}} - 1$ in line 11, which is correct. If d matches q_{from} in line 5, then q_{best} gets the value q_{from} and we enter the if-test on line 7. We need to consider two cases:

(1) *q_{from} is a last sibling:* In this case, we return q_{from} on line 9. This is correct, as $q_{\text{from}} + 1$ is q_{from} 's parent, which cannot be matched onto d due to the semantics of the descendant axis.

(2) *q_{from} is not a last sibling:* If q_{from} has a right sibling, we execute TMATCH recursively on d , $q_{\text{from}} + 1$, and $\text{lastSib}(q_{\text{from}})$, yielding q . By induction, q is computed correctly. That is, if $q = (q_{\text{from}} + 1) - 1$, which implies that d does not match $q_{\text{from}} + 1$, we return q_{from} , which is correct. Otherwise, we argue that $\text{subtree}(d) = d$ matches $[q_{\text{from}}, q]$ but not $\text{subtree}(q + 1)$. By Observation 7, this would complete the proof. By induction, we immediately have that d matches $[q_{\text{from}}, q]$. If $q < \text{lastSib}(q_{\text{from}})$, we also have by induction that d does not match $\text{subtree}(q + 1)$. If $q = \text{lastSib}(q_{\text{from}})$, then $q + 1$ is q_{from} 's parent. Hence, d does not match $\text{subtree}(q + 1)$, as $q + 1$ has a child and d has not. \square

Lemma 7. *Let d be a data node with height $n > 1$ and q_{from} and q_{until} be query nodes. If HMATCH is correct for all data nodes of height up to $n - 1$, then TMATCH is correct for all data nodes of height up to n . That is, given d , q_{from} , and q_{until} , TMATCH returns the largest node q in $[q_{\text{from}}, q_{\text{until}}]$ such that $\text{subtree}(d) \models [q_{\text{from}}, q]$ if it exists; and $q_{\text{from}} - 1$ otherwise.*

PROOF. Assume that HMATCH is correct for all data nodes of height up to $n - 1$. As d is not a leaf, we start by calling HMATCH on $\text{lastChild}(d)$, q_{from} , and q_{until} on line 3 (see also Figure 4(a)), yielding q_{best} . By our assumption, q_{best} is computed correctly. We now prove the lemma by induction on the number of nodes of $[q_{\text{from}}, q_{\text{until}}]$.

$q_{\text{from}} = q_{\text{until}}$: We consider two cases.

(1) If $\text{subhedg}(\text{lastChild}(d))$ does not match q_{from} , then q_{best} is $q_{\text{from}} - 1$. Consequently, we test whether d matches q_{from} on line 5. If d does not match q_{from} , we return $q_{\text{from}} - 1$ on line 11. If d matches q_{from} , then q_{best} gets the value q_{from} . As

$q_{\text{from}} = q_{\text{until}}$ is a last sibling (Observation 5), we do not execute the recursive call on line 8 and return q_{from} in line 9. Both cases are easily seen to be correct.

(2) Otherwise, $q_{\text{best}} = q_{\text{from}} = q_{\text{until}}$. In this case we return q_{best} , which is correct.

$q_{\text{from}} < q_{\text{until}}$: (1) If both $\text{subhedge}(\text{lastChild}(d))$ and d do not match q_{from} , then we return $q_{\text{from}} - 1$ on line 11, which is correct.

(2) If $\text{subhedge}(\text{lastChild}(d))$ matches q_{from} and $q_{\text{best}} = q_{\text{until}}$ on line 5, then we return q_{until} . Due to the correctness of HMATCH, this means that $\text{subhedge}(\text{lastChild}(d))$ already matches $[q_{\text{from}}, q_{\text{until}}]$, hence, $\text{subtree}(d)$ matches $[q_{\text{from}}, q_{\text{until}}]$ by our tree pattern matching semantics.

(3) If $\text{subhedge}(\text{lastChild}(d))$ matches q_{from} , $q_{\text{best}} + 1 \leq q_{\text{until}}$, and d does not match $q_{\text{best}} + 1$ on line 5, then we return q_{best} in line 11. We consider two cases.

- *q_{best} is not a last sibling:* Hence, $q_{\text{best}} + 1$ is a leaf (Observation 4). Due to the correctness of HMATCH for $\text{subhedge}(\text{lastChild}(d))$, we know that $\text{subhedge}(\text{lastChild}(d))$ does not match $\text{subtree}(q_{\text{best}} + 1) = q_{\text{best}} + 1$. Hence, returning q_{best} is correct.
- *q_{best} is a last sibling:* Hence, $q_{\text{best}} + 1$ is q_{best} 's parent. Due to the correctness of HMATCH, we have that $\text{subhedge}(\text{lastChild}(d)) \models [q_{\text{from}}, q_{\text{best}}]$. Towards a contradiction, assume that $\text{subhedge}(d) \models \text{subtree}(q_{\text{best}} + 1)$. As d does not match $q_{\text{best}} + 1$, this implies that $\text{subhedge}(\text{lastChild}(d)) \models \text{subtree}(q_{\text{best}} + 1)$. However, this contradicts that HMATCH is correct. Hence, it is correct to return q_{best} due to Observation 7.

(4) Otherwise, denote by q_{best}^0 the value of the variable q_{best} after the assignment on line 3. We have that q_{best}^0 is correctly computed on line 3 and that d matches $q_{\text{best}}^0 + 1$, after which q_{best} gets the value $q_{\text{best}}^0 + 1$. Notice that $q_{\text{best}}^0 + 1 \geq q_{\text{from}}$. We need to consider two cases:

- *$q_{\text{best}}^0 + 1$ is a last sibling:* We return $q_{\text{best}}^0 + 1$ in line 9. If $q_{\text{best}}^0 + 1 = q_{\text{until}}$, this is correct. If $q_{\text{best}}^0 + 1 < q_{\text{until}}$, towards a contradiction, assume that $\text{subtree}(d)$ matches $\text{subtree}(q_{\text{best}}^0 + 2)$. As $q_{\text{best}}^0 + 2$ is the parent of $q_{\text{best}}^0 + 1$, this would mean that $\text{subhedge}(\text{lastChild}(d)) \models \text{subtree}(q_{\text{best}}^0 + 1)$, which is a contradiction.
- *$q_{\text{best}}^0 + 1$ is not a last sibling:* If $q_{\text{best}}^0 + 1$ has a right sibling, we execute TMATCH on d , $q_{\text{best}}^0 + 2$, and $\text{lastSib}(q_{\text{best}}^0 + 1)$ on line 8, yielding q . By induction, q is computed correctly. If q is $(q_{\text{best}}^0 + 2) - 1$, which implies that $\text{subtree}(d)$ does not match $q_{\text{best}}^0 + 2$, we return $q_{\text{best}}^0 + 1$, which is correct. Otherwise, according to Observation 7, we need to show that $\text{subtree}(d)$ matches $[q_{\text{from}}, q]$ but not $\text{subtree}(q + 1)$. By induction, we have that $\text{subtree}(d)$ matches $[q_{\text{from}}, q]$. If $q < \text{lastSib}(q_{\text{best}}^0 + 1)$, we also have by induction that $\text{subtree}(d)$ does not match $\text{subtree}(q + 1)$. If $q = \text{lastSib}(q_{\text{best}}^0 + 1)$, we have that $\text{subtree}(d)$ does not match $\text{subtree}(q + 1)$, because there does not exist an $u \neq 1$ s.t. $\text{subtree}(d) \models \text{subtree}(q_{\text{best}}^0 + 1)$, and $q + 1$ is $q_{\text{best}}^0 + 1$'s parent. \square

4.1.2. Correctness of HMatch.

Lemma 8. *Let $r_{\text{top}} = \text{RTOP}(q_1, q_2)$ and $q_1 \leq q_2$. If $q_1 \in \text{subhedge}(q_2)$, then $r_{\text{top}} = q_2$ and $q_2 \leq \text{lastSib}(r_{\text{top}})$. If $q_1 \notin \text{subhedge}(q_2)$, then $r_{\text{top}} < q_2$ and $q_2 < \text{lastSib}(r_{\text{top}})$.*

PROOF. Recall that, by definition, $\text{subhedge}(q_2)$ is the interval $[q_{\text{small}}, q_2]$, where q_{small} is the smallest descendant of q_2 's leftmost sibling.

$q_1 \in \text{subhedge}(q_2)$: As both q_1 and q_2 are in $\text{subhedge}(q_2)$, we have that $[q_1, q_2]$ is entirely contained in $\text{subhedge}(q_2)$.

By definition, r_{top} is the largest node in $[q_1, q_2]$ among the nodes with minimal depth. As q_2 has minimal depth in $\text{subhedge}(q_2)$ and q_2 is the largest node in $[q_1, q_2]$, we have that $r_{\text{top}} = q_2$.

$q_1 \notin \text{subhedge}(q_2)$: Notice that this can only occur when q_2 has a parent. As $q_1 \leq q_2$, we have that $q_1 < q_{\text{small}}$. By definition of the left-to-right postordering, we have that q_1 is either a left sibling of an ancestor of q_2 (not including the ancestors themselves), or a descendant-or-self thereof. Let u_1 and u_2 be the two unique siblings such that $u_1 \neq u_2$, q_1 is in $\text{subtree}(u_1)$, and q_2 is in $\text{subtree}(u_2)$. Notice that $q_1 \leq u_1 < q_2 < u_2$. Hence, u_1 is in $[q_1, q_2]$ and $\text{depth}(u_1) < \text{depth}(q_2)$. As q_2 has minimal depth in $\text{subhedge}(q_2)$, we have that r_{top} is not in $\text{subhedge}(q_2)$. By definition of RTOP , this immediately implies that $r_{\text{top}} < q_2$. Furthermore, as $\text{depth}(\text{lastSib}(r_{\text{top}})) = \text{depth}(r_{\text{top}}) \leq \text{depth}(u_1) = \text{depth}(u_2)$ and as $\text{lastSib}(r_{\text{top}})$ is also r_{top} 's largest sibling, we have that $\text{lastSib}(r_{\text{top}}) \geq u_2 > q_2$. \square

Corollary 1. *If $r_{\text{top}} = \text{RTOP}(q_1, q_2)$ then $q_1 \in \text{subhedge}(r_{\text{top}})$.*

PROOF. As r_{top} is in $[q_1, q_2]$, r_{top} is also the rightmost node among the topmost nodes in $[q_1, r_{\text{top}}]$. If we assume that $q_1 \notin \text{subhedge}(r_{\text{top}})$, then Lemma 8 implies that $r_{\text{top}} < r_{\text{top}}$ which is a contradiction. \square

Lemma 9. *All function calls of $\text{TMATCH}(d, q_1, q_2)$ in the loop of HMATCH have the property that $[q_1, q_2]$ is an interval which includes $\text{subtree}(q_{\text{hedge}} + 1)$. All function calls of $\text{HMATCH}(d, q_1, q_2)$ in the loop of HMATCH have the property that $[q_1, q_2]$ is an interval which includes $\text{subtree}(q_{\text{tree}} + 1)$.*

PROOF. For the first statement, we have to show that (i) $q_1 \leq q_{\text{small}}$, where q_{small} is the smallest node in $\text{subtree}(q_{\text{hedge}} + 1)$ and (ii) $q_2 \geq q_{\text{hedge}} + 1$.

First, observe that the function calls of RTOP on line 9 and line 12 results in a value of r_{top} that is at most q_{hedge} . If $r_{\text{top}} < q_{\text{hedge}}$ then $r_{\text{top}} < q_{\text{small}}$ as, by Lemma 8, $r_{\text{top}} = q_{\text{hedge}}$ when r_{top} is in $[q_{\text{small}}, q_{\text{hedge}}]$. Hence, $r_{\text{top}} + 1 \leq q_{\text{small}}$. If $r_{\text{top}} = q_{\text{hedge}}$, we know that r_{top} is not a last sibling due to the condition of the while loop on line 10. Hence, $q_{\text{small}} = q_{\text{hedge}} + 1 = r_{\text{top}} + 1$ is a leaf (Observation 4). This proves property (i).

Property (ii) is immediate as the condition of the while-loop on line 10 requires that $q_2 = \text{lastSib}(r_{\text{top}}) \geq q_{\text{hedge}} + 1$.

The proof of the second statement is analogous to the proof of the first statement. \square

Lemma 10. *The loop on line 6, and the while loops on lines 10 and 18 perform at most a linear number of iterations.*

PROOF. Notice that we exit the loop on line 6 if $\max(q_{\text{tree}}, q_{\text{hedge}})$ does *not* increase. However, this value cannot keep increasing indefinitely as it is bounded from above by q_{until} in the algorithm. Hence, the loop performs at most a linear number of iterations.

The while loop on line 10 terminates after a linear number of iterations, as the value of r_{top} increases with each execution and the while loop only continues as long as r_{top} is smaller than q_{hedge} , a value which remains unchanged. The argument for the while loop on line 18 is analogous. \square

Lemma 11. *Let d be a data node and q_{from} and q_{until} be query nodes. If TMATCH is correct for all data nodes of height up to n , then HMATCH is correct for all data nodes of height up to n . That is, given d , q_{from} , and q_{until} , HMATCH returns the largest node q in $[q_{\text{from}}, q_{\text{until}}]$ such that $\text{subhedge}(d)$ matches $[q_{\text{from}}, q]$ if it exists; and nil otherwise.*

PROOF. Let k be such that d has k left siblings (including d itself). We prove the lemma by induction on k .

If $k = 1$ then the Lemma is immediate from the function call on line 2 and the assumption that TMATCH is correct for all data nodes of height up to n .

So, from now on, we assume that $k > 1$. We need to show that the algorithm returns $q_{\text{from}} - 1$ if $\text{subhedge}(d)$ does not match q_{from} . Otherwise, we show that we return a q in $[q_{\text{from}}, q_{\text{until}}]$ if $\text{subhedge}(d)$ matches $[q_{\text{from}}, q]$ and either

- $q = q_{\text{until}}$, or
- neither $\text{subtree}(d)$, nor $\text{subhedge}(\text{prevSib}(d))$ matches $\text{subtree}(q + 1)$.

In the remainder of the proof, we refer to the above property with the label (\dagger) . The correctness of property (\dagger) follows directly from our tree pattern query semantics if we return $q_{\text{from}} - 1$ and from Observation 7 otherwise. Indeed, from Observation 5 we know that $[q_{\text{from}}, q_{\text{until}}]$ is complete. Furthermore, $\text{subhedge}(d)$ does not match $\text{subtree}(q + 1)$ if and only if neither $\text{subtree}(d)$ nor $\text{subhedge}(\text{prevSib}(d))$ match $\text{subtree}(q + 1)$.

Notice that the loop on line 6 terminates by Lemma 10. We now proceed with an induction over the number ℓ of loop executions proving that the following invariants hold:

- (I1):** if q_{tree} is not $q_{\text{from}} - 1$ then $\text{subhedge}(d)$ matches $[q_{\text{from}}, q_{\text{tree}}]$;
- (I2):** if q_{hedge} is not $q_{\text{from}} - 1$ then $\text{subhedge}(d)$ matches $[q_{\text{from}}, q_{\text{hedge}}]$;
- (I3):** $q_{\text{tree}} = q_{\text{until}}$ or $\text{subtree}(d)$ does not match $\text{subtree}(q_{\text{tree}} + 1)$; and,
- (I4):** $q_{\text{hedge}} = q_{\text{until}}$ or $\text{subhedge}(\text{prevSib}(d))$ does not match $\text{subtree}(q_{\text{hedge}} + 1)$.

At the same time, we show that, if the algorithm returns a certain value q , the property (\dagger) holds for q .

$\ell = 0$ (before the first loop execution): We computed q_{hedge} , which results from executing HMATCH on $\text{prevSib}(d)$, q_{from} , and q_{until} ; and we computed q_{tree} , which results from executing TMATCH on d , q_{from} , and q_{until} (see also Figure 4(c)). By induction on k , we have that q_{hedge} is computed correctly. Moreover, as we assume that TMATCH is correct for all data nodes of height up to n , we also have that q_{tree} is computed correctly. Properties (I1)–(I2) immediately follow from the correctness of the recursive calls of TMATCH and HMATCH. Moreover, Observation 7 implies that (I3)–(I4) also hold. As

the algorithm does not return anything up to here, we do not have to show yet that (\dagger) holds.

$\ell \geq 1$ (*subsequent loop executions*): We consider three cases.

(1) If $q_{\text{hedge}} = q_{\text{tree}}$, we return q_{hedge} . This is correct, as in this case, properties (I1)–(I4) immediately imply property (\dagger) .

(2) If $q_{\text{tree}} < q_{\text{hedge}}$, notice that we do not change the value of q_{hedge} in this iteration of the loop. Hence, for the induction, we only need to show that properties (I1) and (I3) are preserved. We consider two cases.

If $q_{\text{hedge}} = q_{\text{until}}$ the while loop in line 10 is not executed and we return q_{until} in line 14. Here, it follows immediately from (I2) that (\dagger) holds.

If $q_{\text{hedge}} < q_{\text{until}}$ we consider two cases.

- If $\text{subtree}(d)$ does not match $\text{subtree}(q_{\text{hedge}}+1)$, none of the function calls $\text{TMATCH}(d, q_1, q_2)$ in the while loop yield a value greater than q_{hedge} . This follows from the correctness of TMATCH for data nodes up to height n , and from Lemma 9, stating that $[q_1, q_2]$ always includes $\text{subtree}(q_{\text{hedge}} + 1)$. Indeed, should such a function call $\text{TMATCH}(d, q_1, q_2)$ yield a greater value than q_{hedge} , then we would have that $\text{subtree}(d)$ matches $\text{subtree}(q_{\text{hedge}} + 1)$, which contradicts that we are investigating the case that $\text{subtree}(d)$ does not match $\text{subtree}(q_{\text{hedge}} + 1)$. Hence, we return q_{hedge} in line 14. Correctness of the property (\dagger) for q_{hedge} now follows from the following facts:

- $q_{\text{hedge}} \geq q_{\text{from}}$, as $q_{\text{tree}} < q_{\text{hedge}}$;
- $q_{\text{hedge}} < q_{\text{until}}$;
- $\text{subhedge}(d)$ matches $[q_{\text{from}}, q_{\text{hedge}}]$, by (I2);
- $\text{subtree}(d)$ does not match $\text{subtree}(q_{\text{hedge}} + 1)$; and,
- $\text{subhedge}(\text{prevSib}(d))$ does not match $\text{subtree}(q_{\text{hedge}} + 1)$ by (I4).

- If $\text{subtree}(d)$ matches $\text{subtree}(q_{\text{hedge}} + 1)$ the proof is more complicated. First, observe that the while loop on line 10 terminates by Lemma 10.

For the remainder of this case, we will show that $q_{\text{tree}} > q_{\text{hedge}}$ after exiting the while loop in the $i + 1^{\text{th}}$ execution of the test on line 10. In particular, this implies that the algorithm will not return any value in iteration ℓ of the loop. So we only need to show that, at the end of the current iteration, properties (I1) and (I3) hold.

To show (I3), we will show that, if in the j^{th} execution of the while loop we obtain a value q for the variable q_{tree} for which it holds that $q > q_{\text{hedge}}$ then we either have that $q = q_{\text{until}}$ or that $\text{subtree}(d)$ does not match $\text{subtree}(q + 1)$. Afterwards, we show (I1).

We start by showing that $q_{\text{tree}} > q_{\text{hedge}}$ after exiting the while loop:

Goal 1: $q_{\text{tree}} > q_{\text{hedge}}$ after exiting the while loop in the $i + 1^{\text{th}}$ execution of the test on line 10. So we execute the while body i times and then exit the loop.

Let q_{tree}^i denote the value of q_{tree} at the end of the i^{th} execution (i.e., after the assignment on line 11) and let q_{tree}^0 be the value of q_{tree} before entering the while loop. Furthermore, let rtop^i denote the value of rtop at the end of the i^{th} execution

(i.e., after the assignment on line 12). Let rtop^0 be the value of rtop before entering the while loop.

($i = 0$): We will show that this case does not occur. That is, the body of the while loop is always executed at least once. Towards a contradiction, assume that we do not execute the body of the while loop. We consider two cases. If we exit the while loop one of them must hold.

- *Case 1:* $\text{rtop}^0 < \infty$ and $q_{\text{hedge}} \geq \text{lastSib}(\text{rtop}^0)$. Recall that $\text{rtop}^0 = \text{RTOP}(q_{\text{tree}}^0 + 1, q_{\text{hedge}})$. Due to Lemma 8, $q_{\text{hedge}} \geq \text{lastSib}(\text{rtop}^0)$ implies that (i) $\text{rtop}^0 = \text{lastSib}(\text{rtop}^0) = q_{\text{hedge}}$ and that (ii) $q_{\text{tree}}^0 + 1$ is in $\text{subhedge}(q_{\text{hedge}})$. As $q_{\text{hedge}} < q_{\text{until}}$ and q_{hedge} is a last sibling this means that $q_{\text{tree}}^0 + 1$ is in $\text{subtree}(q_{\text{hedge}} + 1)$. Moreover, as we are in the case that $q_{\text{tree}} < q_{\text{hedge}}$, we know by induction on ℓ (statement (I3) in particular) that $\text{subtree}(d)$ does not match $\text{subtree}(q_{\text{tree}}^0 + 1)$. However, as we have shown above that $q_{\text{tree}}^0 + 1$ is in $\text{subtree}(q_{\text{hedge}} + 1)$, this contradicts the fact that we are in the case that $\text{subtree}(d)$ matches $\text{subtree}(q_{\text{hedge}} + 1)$.
- *Case 2:* $\text{rtop}^0 = \infty$. By definition of RTOP , this means that $q_{\text{tree}}^0 + 1 > q_{\text{hedge}}$. But we are currently investigating in the case that $q_{\text{tree}}^0 < q_{\text{hedge}}$. Contradiction.

Hence, we showed that the while loop on line 10 is executed at least once.

($i > 0$): Again, we consider the two possible settings in which we exit the while loop. We show again that the first of the two does not occur here.

- *Case 1:* $\text{rtop}^i < \infty$ and $q_{\text{hedge}} \geq \text{lastSib}(\text{rtop}^i)$. Recall that $\text{rtop}^i = \text{RTOP}(q_{\text{tree}}^i + 1, q_{\text{hedge}})$. Due to Lemma 8, $q_{\text{hedge}} \geq \text{lastSib}(\text{rtop}^i)$ implies that (i) $\text{rtop}^i = \text{lastSib}(\text{rtop}^i) = q_{\text{hedge}}$ and that (ii) $q_{\text{tree}}^i + 1$ is in $\text{subhedge}(q_{\text{hedge}})$, implying that $q_{\text{tree}}^i + 1 \leq q_{\text{hedge}}$. As $q_{\text{hedge}} < q_{\text{until}}$ and q_{hedge} is a last sibling this means that $q_{\text{tree}}^i + 1$ is in $\text{subtree}(q_{\text{hedge}} + 1)$. Since we did not exit the while loop in the i th test, we have that $q_{\text{hedge}} < \text{lastSib}(\text{rtop}^{i-1})$. Hence, we have that $q_{\text{tree}}^i + 1 \leq q_{\text{hedge}} < \text{lastSib}(\text{rtop}^{i-1})$. Recall that $q_{\text{tree}}^i = \text{HMATCH}(\text{prevSib}(d), \text{rtop}^{i-1} + 1, \text{lastSib}(\text{rtop}^{i-1}))$. By the correctness of TMATCH , Observation 7, and the fact that $[\text{rtop}^{i-1} + 1, \text{lastSib}(\text{rtop}^{i-1})]$ is a complete interval (Observation 5) we can conclude that $\text{subtree}(d)$ does not match $\text{subtree}(q_{\text{tree}}^i + 1)$ which, we argued above, is a subtree of $\text{subtree}(q_{\text{hedge}} + 1)$. Hence, $\text{subtree}(d)$ does not match $\text{subtree}(q_{\text{hedge}} + 1)$, which contradicts the fact that we are in the case that $\text{subtree}(d)$ matches $\text{subtree}(q_{\text{hedge}} + 1)$.
- *Case 2:* $\text{rtop}^i = \infty$. Hence, $q_{\text{tree}}^i + 1 > q_{\text{hedge}}$. We prove that it cannot be the case that $q_{\text{tree}}^i = q_{\text{hedge}}$. Hence, $q_{\text{tree}}^i > q_{\text{hedge}}$ and Goal 1 follows. To this end, assume, towards a contradiction, that $q_{\text{tree}}^i = q_{\text{hedge}}$. Recall that $q_{\text{tree}}^i = \text{TMATCH}(d, \text{rtop}^{i-1} + 1, \text{lastSib}(\text{rtop}^{i-1}))$. Moreover, $\text{lastSib}(\text{rtop}^{i-1}) > q_{\text{hedge}}$ since otherwise we would have exited the while loop right after test i . We conclude that $q_{\text{hedge}} + 1$ is a node in $[\text{rtop}^{i-1} + 1, \text{lastSib}(\text{rtop}^{i-1})]$. However, as $\text{subtree}(d)$ matches $\text{subtree}(q_{\text{hedge}} + 1)$, this would imply that $\text{subtree}(d)$ also matches $[\text{rtop}^{i-1} + 1, q_{\text{hedge}} + 1] = [\text{rtop}^{i-1} + 1, q_{\text{tree}}^i + 1]$ which is in contradiction with the correctness of TMATCH .

This concludes the proof of Goal 1.

Goal 2. *If in the j^{th} execution of the while loop we obtain a value q for the variable q_{tree} for which it holds that $q > q_{\text{hedge}}$ and $q+1 \leq q_{\text{until}}$, then we have that $\text{subtree}(d)$ does not match $\text{subtree}(q+1)$.*

Observe that we need at least one execution of the body of the while, since before the first execution we have that $q_{\text{tree}} < q_{\text{hedge}}$. Let q_{tree}^j denote the value of q_{tree} at the end of the j^{th} execution (i.e., after the assignment on line 11) and let q_{tree}^0 be the value of q_{tree} before entering the while loop. Furthermore let rtop^j denote the value of rtop at the end of the j^{th} execution (i.e., after the assignment on line 12). Let rtop^0 be the value of rtop before entering the while loop.

Hence, for every $j \geq 1$, q_{tree}^j is the result of a function call $\text{TMATCH}(d, \text{rtop}^{j-1} + 1, \text{lastSib}(\text{rtop}^{j-1}))$. If $q_{\text{tree}}^j > q_{\text{hedge}}$ we will exit the while loop right after the current iteration. We consider three cases.

- If $q_{\text{tree}}^j < \text{lastSib}(\text{rtop}^{j-1})$ we have that $\text{subtree}(d)$ does not match the subtree of $q_{\text{tree}}^j + 1$ due to the correctness of TMATCH for data nodes up to height n and Observation 7.
- If $q_{\text{tree}}^j = q_{\text{until}}$ the claim is trivial.
- The remaining case is that $q_{\text{tree}}^j = \text{lastSib}(\text{rtop}^{j-1}) < q_{\text{until}}$. In this case, $q_{\text{tree}}^j + 1$ is the parent of q_{tree}^j due to Observation 4. We consider two cases.
 $j = 1$: We want to prove that $\text{subtree}(d)$ does not match $\text{subtree}(q_{\text{tree}}^0 + 1)$ and that $\text{subtree}(q_{\text{tree}}^0 + 1)$ is a subtree of $\text{subtree}(q_{\text{tree}}^1 + 1)$. Then we can conclude that $\text{subtree}(d)$ does not match $\text{subtree}(q_{\text{tree}}^1 + 1)$.

We start by proving that $\text{subtree}(d)$ does not match $\text{subtree}(q_{\text{tree}}^0 + 1)$. By induction on ℓ (and, in particular, by (I3)) we know that $q_{\text{tree}}^0 = q_{\text{until}}$ or $\text{subtree}(d)$ does not match $\text{subtree}(q_{\text{tree}}^0 + 1)$. If $q_{\text{tree}}^0 = q_{\text{until}}$ we wouldn't be in the case that $q_{\text{tree}}^0 < q_{\text{hedge}}$. We can conclude that $\text{subtree}(d)$ does not match $\text{subtree}(q_{\text{tree}}^0 + 1)$.

It remains to be shown that $\text{subtree}(q_{\text{tree}}^0 + 1)$ is a subtree of $\text{subtree}(q_{\text{tree}}^1 + 1)$. Line 9 states that $\text{rtop}^0 = \text{RTOP}(q_{\text{tree}}^0 + 1, q_{\text{hedge}})$. Corollary 1 implies that then $q_{\text{tree}}^0 + 1$ is a node in $\text{subhedge}(\text{rtop}^0)$. Now we take into consideration that we are investigating in the case that $q_{\text{tree}}^1 = \text{lastSib}(\text{rtop}^0)$ which implies that $\text{subhedge}(\text{rtop}^0) \subseteq \text{subhedge}(q_{\text{tree}}^1)$. Combining this with the consequence of the Corollary it follows that $q_{\text{tree}}^0 + 1$ is a node in $\text{subhedge}(q_{\text{tree}}^1)$. Recall that $q_{\text{tree}}^1 + 1$ is q_{tree}^1 's parent. Hence, $q_{\text{tree}}^0 + 1$ is a node in $\text{subtree}(q_{\text{tree}}^1 + 1)$ and $\text{subtree}(q_{\text{tree}}^0 + 1)$ is a subtree of $\text{subtree}(q_{\text{tree}}^1 + 1)$.

$j > 1$: Analogously as in the $j = 1$ case, we prove that $\text{subtree}(d)$ does not match $\text{subtree}(q_{\text{tree}}^{j-1} + 1)$ and that $\text{subtree}(q_{\text{tree}}^{j-1} + 1)$ is a subtree of $\text{subtree}(q_{\text{tree}}^j + 1)$. Then we can conclude that $\text{subtree}(d)$ does not match $\text{subtree}(q_{\text{tree}}^j + 1)$.

We start by proving that $\text{subtree}(d)$ does not match $\text{subtree}(q_{\text{tree}}^{j-1} + 1)$. We have that $q_{\text{tree}}^{j-1} = \text{TMATCH}(d, \text{rtop}^{j-2} + 1, \text{lastSib}(\text{rtop}^{j-2}))$. Notice that, if $q_{\text{tree}}^{j-1} < \text{lastSib}(\text{rtop}^{j-2})$, we immediately have by the correctness of TMATCH and Observation 7 that $\text{subtree}(d)$ does not match $\text{subtree}(q_{\text{tree}}^{j-1} + 1)$. So, towards a contradiction, let us assume that $q_{\text{tree}}^{j-1} \geq \text{lastSib}(\text{rtop}^{j-2})$.

Notice that $q_{\text{hedge}} < \text{lastSib}(\text{rtop}^{j-2})$ and that $\text{rtop}^{j-1} \leq q_{\text{hedge}}$, otherwise we wouldn't have arrived in the j th iteration. Moreover, $\text{rtop}^{j-2} < \text{rtop}^{j-1}$. As $\text{rtop}^{j-2} < \text{rtop}^{j-1} \leq \text{lastSib}(\text{rtop}^{j-2})$, we also have that $\text{lastSib}(\text{rtop}^{j-1}) \leq \text{lastSib}(\text{rtop}^{j-2}) < q_{\text{tree}}^{j-1} + 1 \leq q_{\text{tree}}^j$, which is in contradiction with $q_{\text{tree}}^j = \text{lastSib}(\text{rtop}^{j-1})$, which is the case we are investigating.

It remains to be shown that $\text{subtree}(q_{\text{tree}}^{j-1} + 1)$ is a subtree of $\text{subtree}(q_{\text{tree}}^j + 1)$. Line 12 states that $\text{rtop}^{j-1} = \text{RTOP}(q_{\text{tree}}^{j-1} + 1, q_{\text{hedge}})$. Corollary 1 implies that then $q_{\text{tree}}^{j-1} + 1$ is a node in $\text{subhedge}(\text{rtop}^{j-1})$. Now we take into consideration that we are investigating the case that $q_{\text{tree}}^j = \text{lastSib}(\text{rtop}^{j-1})$ which implies that $\text{subhedge}(\text{rtop}^{j-1}) \subseteq \text{subhedge}(q_{\text{tree}}^j)$. Combining this with the consequence of the Corollary it follows that $q_{\text{tree}}^{j-1} + 1$ is a node in $\text{subhedge}(q_{\text{tree}}^j)$. Recall that $q_{\text{tree}}^j + 1$ is q_{tree}^j 's parent. Hence, $q_{\text{tree}}^{j-1} + 1$ is a node in $\text{subtree}(q_{\text{tree}}^j + 1)$ and $\text{subtree}(q_{\text{tree}}^{j-1} + 1)$ is a subtree of $\text{subtree}(q_{\text{tree}}^j + 1)$.

This concludes the proof of Goal 2.

It remains to show that (I1) holds at the end of the ℓ -th iteration of the loop, that is, that $\text{subhedge}(d)$ matches $[q_{\text{from}}, q_{\text{tree}}]$. Due to (I2) we have that $\text{subhedge}(d)$ matches $[q_{\text{from}}, q_{\text{hedge}}]$. Recall that the number of while loop executions is at least one. Hence, we have that $q_{\text{tree}} = \text{TMATCH}(d, \text{rtop} + 1, \text{lastSib}(\text{rtop}))$, where $\text{rtop} \leq q_{\text{hedge}} < q_{\text{tree}} \leq \text{lastSib}(\text{rtop})$. The first inequality follows from the fact that $\text{rtop} < \infty$ and the definition of RTOP , the second one follows from Goal 1, and the third one from the correctness of TMATCH . Hence, we have that

- $\text{subhedge}(d) \models [q_{\text{from}}, \text{rtop}]$ and
- $\text{subtree}(d) \models [\text{rtop} + 1, q_{\text{tree}}]$.

Moreover, the facts that $\text{rtop} + 1$ is a leaf (Observation 4) and $q_{\text{tree}} \leq \text{lastSib}(\text{rtop})$ imply that $\text{subhedge}(d) \models [q_{\text{from}}, q_{\text{tree}}]$.

This concludes the proof the case where $\text{subtree}(d)$ matches $\text{subtree}(q_{\text{hedge}} + 1)$.

This concludes the proof of the case where $q_{\text{hedge}} < q_{\text{until}}$, and also the proof of the case where $q_{\text{tree}} < q_{\text{hedge}}$.

(3) If $q_{\text{hedge}} < q_{\text{tree}}$ the proof is dual to the proof of case (2). \square

The correctness of Lemma 5 now follows from Lemmas 6, 7, and 11.

We now argue how TMATCH can be modified to a procedure TMATCH-ALL , that computes *all* data nodes u such that $D \models^u Q$. In order to compute *all* the matches, we add a test to line 9 of TMATCH . That is, before returning q_{best} , we test whether q_{best} is the root of Q , and we output d if it is. Now we return $q_{\text{best}} - 1$, as if the query root was not matched. Furthermore, TMATCH-ALL recursively calls TMATCH-ALL and HMATCH-ALL instead of TMATCH and HMATCH . Here HMATCH-ALL is the same as HMATCH , except that it recursively calls TMATCH-ALL and HMATCH-ALL instead of HMATCH and TMATCH .

The following theorem can now be proved:

Theorem 8. *Let d be the root node of D and let q_{from} be the smallest and q_{root} be the largest node of Q , respectively. TMATCH-ALL is correct, that is, $\text{TMATCH-ALL}(d, q_{\text{from}}, q_{\text{until}})$ outputs the data nodes u such that $D \models^u Q$.*

PROOF. It follows directly from our additional test and the correctness of TMATCH that $D \models^u Q$ for all the nodes u that TMATCH-ALL outputs.

It remains to prove that, if $D \models^u Q$, then TMATCH-ALL outputs u . Towards a contradiction, assume that there is an u such that $D \models^u Q$, but u was not reported by TMATCH-ALL. By an easy induction it can be shown that for every data node d_0 in D there is a call TMATCH-ALL for d_0 's subtree and Q . In particular, there was a call TMATCH-ALL($u, q_{\text{from}}, q_{\text{root}}$). Since this call did not output u , it follows that u must have children and that $\text{HMATCH-ALL}(\text{lastChild}(u), q_{\text{from}}, q_{\text{root}}) < q_{\text{root}} - 1$, (because otherwise q_{root} and u would have been compared and u would have been written to the output). In general, we have that $\text{HMATCH-ALL}(d, q_1, q_2) = \min((\text{HMATCH}(d, q_1, q_2), q_{\text{root}} - 1))$. It then follows that $\text{HMATCH-ALL}(\text{lastChild}(u), q_{\text{from}}, q_{\text{root}}) = \text{HMATCH}(\text{lastChild}(u), q_{\text{from}}, q_{\text{root}})$.

If we now call TMATCH($u, q_{\text{from}}, q_{\text{root}}$), it calls $\text{HMATCH}(\text{lastChild}(u), q_{\text{from}}, q_{\text{root}})$, which yields again a value less than $q_{\text{root}} - 1$. Therefore, the return value of TMATCH($u, q_{\text{from}}, q_{\text{root}}$) is less than q_{root} . But we assumed that $\text{subtree}(u) \models Q$, which contradicts the correctness of TMATCH proved in Lemma 5. \square

4.2. Time and Space Complexity.

First, we need to show that our algorithm determines in PTIME whether $D \models Q$. Notice that the naïve manner of computing the running time of TMATCH gives rise to only an exponential upper bound. Indeed, define (i) $T(N)$ as the running time of TMATCH on d, q_{from} , and q_{until} , where $\text{subtree}(d)$ and $[q_{\text{from}}, q_{\text{until}}]$ have N nodes in total, and (ii) $H(N)$ as the running time of HMATCH on d, q_{from} , and q_{until} , where $\text{subhedger}(d)$ and $[q_{\text{from}}, q_{\text{until}}]$ have N nodes in total. Then, we have that $T(2) \leq p(N)$ for a polynomial p , $T(N) \leq p(N) + H(N - 1) + T(N - 1)$, and $H(N) \leq T(N) + X(N)$, where $X(N) \geq 0$. Hence, $T(N) \leq 2^{N-1}$, which is obviously not sufficient.

We therefore employ a slightly more sophisticated approach in the following Lemma.

Lemma 12. *Given the root node of a data tree D , and the smallest and largest query nodes and of a query tree Q , respectively, TMATCH runs in time $\mathcal{O}(|D| \cdot |Q| \cdot \text{depth}(Q))$. Moreover, TMATCH makes $\mathcal{O}(|D| \cdot |Q|)$ comparisons between a data node and a query node.*

PROOF. Let $|D|$ and $|Q|$ be the number of nodes in the data and query tree, respectively. We first show by induction on the height n of the data node d that the number of calls to the function TMATCH in the computation tree is at most $|D||Q|$. To this end, we prove three intermediate goals.

Goal 1: *Let d be a leaf data node. A computation of TMATCH($d, q_{\text{from}}, q_{\text{until}}$) yielding result q makes at most $\lceil [q_{\text{from}}, q + 1] \rceil$ calls to TMATCH.*

By induction on the size of the query tree interval $[q_{\text{from}}, q_{\text{until}}]$. If d is a leaf and $q_{\text{from}} = q_{\text{until}}$, then TMATCH does not call HMATCH recursively and the test on line 7 fails. Therefore, there is only 1 call to TMATCH and the induction hypothesis holds. If $q_{\text{from}} < q_{\text{until}}$, and TMATCH is not called recursively, then the minimal value we return is $q_{\text{from}} - 1$. Again, there is only 1 call to TMATCH and the induction hypothesis holds. Otherwise, we call TMATCH on line 8, yielding result q . By induction, the total number of calls to TMATCH is at most $1 + \lceil [q_{\text{from}} + 1, q + 1] \rceil$. As $\lceil [q_{\text{from}}, q + 1] \rceil = 1 + \lceil [q_{\text{from}} + 1, q + 1] \rceil$, the induction holds. This concludes the proof of Goal 1.

Goal 2: Let d be a data node with height $n > 1$. If the computation of $\text{HMATCH}(\text{lastChild}(d), q_{\text{from}}, q_{\text{until}})$, yielding the result q_{best}^0 , performs at most $|\text{subhedge}(\text{lastChild}(d))| \cdot |[q_{\text{from}}, q_{\text{best}}^0 + 1]|$ calls to TMATCH , then the computation of $\text{TMATCH}(d, q_{\text{from}}, q_{\text{until}})$, yielding result q , makes at most $|\text{subtree}(d)| \cdot |[q_{\text{from}}, q + 1]|$ calls to TMATCH .

We prove Goal 2 by induction on the size of the query tree interval $[q_{\text{from}}, q_{\text{until}}]$. TMATCH starts by calling $\text{HMATCH}(\text{lastChild}(d), q_{\text{from}}, q_{\text{until}})$ yielding q_{best}^0 . Hence, $|\text{subhedge}(\text{lastChild}(d))| \cdot |[q_{\text{from}}, q_{\text{best}}^0 + 1]|$ calls to TMATCH are performed by this sub-routine.

If $q_{\text{from}} = q_{\text{until}}$, then we either return q_{best}^0 on line 11 or $q_{\text{best}}^0 + 1$ on line 9. In both cases, the number of calls to TMATCH is at most $|\text{subhedge}(\text{lastChild}(d))| \cdot |[q_{\text{from}}, q_{\text{best}}^0 + 1]| + 1$ which is at most $|\text{subtree}(d)| \cdot |[q_{\text{from}}, q_{\text{best}}^0 + 1]|$.

If $q_{\text{from}} < q_{\text{until}}$, and TMATCH is not called recursively, then the minimal value we return is q_{best}^0 . Again, the number of calls to TMATCH is at most $1 + |\text{subhedge}(\text{lastChild}(d))| \cdot |[q_{\text{from}}, q_{\text{best}}^0 + 1]|$ and the induction hypothesis holds. Otherwise, we call TMATCH on line 8, yielding result q . By induction, the total number of calls to TMATCH is at most $1 + |\text{subhedge}(\text{lastChild}(d))| \cdot |[q_{\text{from}}, q_{\text{best}}^0 + 1]| + |\text{subtree}(d)| \cdot |[q_{\text{best}}^0 + 2, q + 1]|$ which is at most $|\text{subtree}(d)| \cdot |[q_{\text{from}}, q + 1]|$. This concludes the proof of Goal 2.

Goal 3: Let d be a data node. If the computation of $\text{TMATCH}(d, q_1, q_2)$, yielding q_{tree} makes at most $|\text{subtree}(d)| \cdot |[q_1, q_{\text{tree}} + 1]|$ calls to TMATCH , then the computation of $\text{HMATCH}(d, q_{\text{from}}, q_{\text{until}})$, yielding q makes at most $|\text{subhedge}(d)| \cdot |[q_{\text{from}}, q + 1]|$ calls to TMATCH .

Let k be such that d has k left siblings (including d itself). We prove the lemma by induction on k . If $k = 1$, Goal 3 is an immediate consequence from the assumption of Goal 3 and the recursive call of TMATCH on line 2. If $k > 1$, then we start by calling $\text{HMATCH}(\text{prevSib}(d), q_{\text{from}}, q_{\text{until}})$, yielding $q_{\text{hedge}}^{1,0}$, and calling $\text{TMATCH}(d, q_{\text{from}}, q_{\text{until}})$, yielding $q_{\text{tree}}^{1,0}$. By induction on k , we have that the call of HMATCH induces $|\text{subhedge}(\text{prevSib}(d))| \cdot |[q_{\text{from}}, q_{\text{hedge}}^{1,0} + 1]|$ calls to TMATCH . Moreover, by the statement of Goal 3, we have that the recursive call of TMATCH induces $|\text{subtree}(d)| \cdot |[q_{\text{from}}, q_{\text{tree}}^{1,0} + 1]|$ calls to TMATCH in total.

According to Lemma 10, the loops on line 6, 10, and 18 perform at most a linear number of iterations. Hence, TMATCH and HMATCH are called (directly) at most a quadratic number of times in the loop.

By $q_{\text{tree}}^{i,j}$, we denote the value of the variable q_{tree} in the i -th iteration of the loop and at the end of the j -th iteration of the while loop in line 10. Moreover, let ℓ denote the number of loop executions and let \max_i denote the number of executions of the while loop on line 10 in the i -th loop execution. Then, we have that every computation of $\text{TMATCH}(d, q_1, q_2)$ in the while loop performs at most $|\text{subtree}(d)| \cdot |[q_{\text{tree}}^{i,j-1} + 2, q_{\text{tree}}^{i,j} + 1]|$ calls to TMATCH when $j > 1$ and at most $|\text{subtree}(d)| \cdot |[q_{\text{tree}}^{i-1, \max_{i-2}(j)} + 2, q_{\text{tree}}^{i,1} + 1]|$ calls otherwise. Notice that $q_{\text{tree}}^{1,0} < q_{\text{tree}}^{1,1} < \dots < q_{\text{tree}}^{1, \max_1} < q_{\text{tree}}^{2,1} < \dots < q_{\text{tree}}^{\ell, \max_\ell} \leq q$, where q is the value we return. Hence, the sum of the calls to TMATCH made by the computations of TMATCH on line 11 is at most $|\text{subtree}(d)| \cdot |[q_{\text{tree}}^{1,0} + 2, q + 1]|$.

Analogously, we obtain that the sum of the calls to TMATCH by the computations of HMATCH on line 19 is at most $|\text{subhedge}(\text{prevSib}(d))| \cdot |[q_{\text{tree}}^{1,0} + 2, q + 1]|$.

In total, this means that the number of calls to `TMATCH` is at most

$$\begin{aligned} & |\text{subhedge}(\text{prevSib}(d))| \cdot |[q_{\text{from}}, q_{\text{hedge}}^{1,0} + 1]| \\ & \quad + |\text{subhedge}(\text{prevSib}(d))| \cdot |[q_{\text{hedge}}^{1,0} + 2, q + 1]| \\ & \quad + |\text{subtree}(d)| \cdot |[q_{\text{from}}, q_{\text{tree}}^{1,0} + 1]| + |\text{subtree}(d)| \cdot |[q_{\text{tree}}^{1,0} + 2, q + 1]| \end{aligned}$$

which is at most $|\text{subhedge}(d)| \cdot |[q_{\text{from}}, q + 1]|$. Hence, Goal 3 follows.

As a consequence of Goals 1, 2, and 3, the total number of calls to `TMATCH` performed by the algorithm is $|D||Q|$. As the only data versus query node comparison in the algorithm occurs in line 5 of `TMATCH`, and as each call of `TMATCH` performs at most one data versus query node comparison (excluding comparisons in recursive calls), the total algorithm also performs at most $|D||Q|$ data versus query node comparisons.

We now argue how this leads us to showing that the overall algorithm has polynomial running time. Consider the entire tree of the calls to `TMATCH` and `HMATCH` in the algorithm, where the children of a node are the functions it calls directly. This computation tree contains at most $|D||Q|$ calls of `TMATCH`. Moreover, every call of `HMATCH` performs at least one direct recursive call to `TMATCH`, so the computation tree also contains at most $|D||Q|$ calls of `HMATCH`. Analogously, the entire computation tree contains at most $|D||Q|$ calls to `rtop`. As `rtop` can be implemented to run in time $\mathcal{O}(\text{depth}(Q))$, the total algorithm runs in time $\mathcal{O}(|D||Q|\text{depth}(Q))$. \square

The $\text{depth}(Q)$ factor in the complexity of `TMATCH` is due to the calls to `rtop` in `HMATCH`, and the computation of the successors of query nodes.

From the complexity of `TMATCH` and the definition of `TMATCH-ALL`, we can immediately infer the following complexity results about `TMATCH-ALL`.

Theorem 9. *`TMATCH-ALL`(D, Q) runs in time $\mathcal{O}(|D| \cdot |Q| \cdot \text{depth}(Q))$. Moreover, `TMATCH-ALL` makes $\mathcal{O}(|D| \cdot |Q|)$ comparisons between a data node and a query node.*

Currently, the maximum recursion depth of `TMATCH-ALL` is $\mathcal{O}(\text{depth}(D) \times \text{branch}(D))$, where $\text{branch}(D)$ is the maximum number of children a node in D has. We have the $\text{branch}(D)$ factor because `HMATCH`($d, q_{\text{from}}, q_{\text{until}}$) calls `HMATCH`($\text{prevSib}(d), q_{\text{from}}, q_{\text{until}}$). However, this bound can be improved using a simple preprocessing step: we can turn D into a binary tree D_{bin} by inserting intermediate levels of special nodes between each data node and its children. By doing so, D only grows linearly in size and the depth only grows by a factor of $\log(\text{branch}(D))$.

As Q only uses descendant axes, we have that $D \models^u Q$ if and only if $D_{\text{bin}} \models^u Q$.⁵ When this preprocessing step is carried out, our algorithm still has $\mathcal{O}(|D||Q|\text{depth}(Q))$ time complexity, but the recursion/stack depth is improved to $\mathcal{O}(\text{depth}(D) \cdot \log(\text{branch}(D)))$.

5. Conclusions and Final Thoughts

As our main results we have exhibited a complexity result, showing that tree pattern matching with only descendant axes is LOGSPACE-complete; and a time-and space

⁵Under the assumption that the new dummy nodes do not match *, which can be trivially incorporated in the algorithm.

efficient bottom-up algorithm for computing all possible exact matches of such a tree pattern in a tree.

From a theory point of view, this is still only a small step in finding the exact complexity of positive conjunctive Core XPath with only child and descendant axes (or, alternatively, tree pattern queries with child and descendant axes), which is probably the most widely used fragment of XPath in practice. Hence, it is quite surprising that the exact complexity of this fragment is still unknown.

From a practical point of view, our bottom-up algorithm gives a good space and time bound on the processing of such descendant-only tree pattern queries. A minor annoyance we still feel for the algorithm is the $\text{depth}(Q)$ factor in the time complexity. However, we need to stress that, in practical applications, $\text{depth}(Q)$ will indeed be very small. In our algorithm, this $\text{depth}(Q)$ factor arises from computing the $\text{RTOP}(q_{\text{tree}}, q_{\text{hedge}})$ -values in each call of HMATCH in the algorithm. It may be possible that this factor can be avoided when integrating the computation of these values in the recursion of the algorithm. For a practical application, one can also avoid the $\text{depth}(Q)$ factor in run-time evaluation by a pre-processing step that computes all the values of $\text{RTOP}(q_{\text{tree}}, q_{\text{hedge}})$ in advance on the query.

References

- [1] M. Altinel and M. Franklin. Efficient filtering of XML documents for selective dissemination of information. In *Proc. International Conference on Very Large Data Bases (VLDB)*, pages 53–64, 2000.
- [2] Z. Bar-Yossef, M. Fontoura, and V. Josifovski. On the memory requirements of XPath evaluation over XML streams. In *Proc. International Symposium on Principles of Database Systems (PODS)*, pages 177–188, 2004.
- [3] Z. Bar-Yossef, M. Fontoura, and V. Josifovski. Buffering in query evaluation over XML streams. In *Proc. International Symposium on Principles of Database Systems (PODS)*, 2005.
- [4] N. Bruno, D. Srivastava, and N. Koudas. Holistic twig joins: Optimal XML pattern matching. In *Proc. ACM SIGMOD International Conference on Management of Data (SIGMOD)*, pages 310–321, 2002.
- [5] C. Y. Chan, W. Fan, P. Felber, M. N. Garofalakis, and R. Rastogi. Tree pattern aggregation for scalable XML data dissemination. In *Proc. International Conference on Very Large Data Bases (VLDB)*, pages 826–837, 2002.
- [6] C. Y. Chan, P. Felber, M. N. Garofalakis, and R. Rastogi. Efficient filtering of XML documents with XPath expressions. In *Proc. International Conference on Data Engineering (ICDE)*, pages 235–244, 2000.
- [7] J. Clark and S. DeRose. XML Path Language (XPath). Technical report, World Wide Web Consortium, November 1999. <http://www.w3.org/TR/xpath>.
- [8] S. A. Cook and P. McKenzie. Problems complete for deterministic logarithmic space. *Journal of Algorithms*, 8:385–394, 1987.
- [9] Y. Diao, M. Altinel, M. J. Franklin, H. Zhang, and P. Fischer. Path sharing and predicate evaluation for high-performance XML filtering. *ACM Transactions on Database Systems*, 28(4):467–516, 2003.
- [10] G. Gottlob, C. Koch, and R. Pichler. Efficient algorithms for processing XPath queries. *ACM Transactions on Database Systems*, 30(2):444–491, 2005.
- [11] G. Gottlob, C. Koch, R. Pichler, and L. Segoufin. The complexity of XPath query evaluation and XML typing. *Journal of the ACM*, 52(2):284–335, 2005.
- [12] G. Gottlob, N. Leone, and F. Scarcello. The complexity of acyclic conjunctive queries. *Journal of the ACM*, 48(1):431–498, 2001.
- [13] M. Götz, C. Koch, and W. Martens. Efficient algorithms for the tree homeomorphism problem. In *Proc. International Symposium on Database Programming Languages (DBPL)*, pages 17–31, 2007.
- [14] T. J. Green, A. Gupta, G. Miklau, M. Onizuka, and D. Suciu. Processing XML streams with deterministic automata and stream indexes. *ACM Transactions on Database Systems*, 29(4):752–788, 2004.

- [15] M. Grohe, C. Koch, and N. Schweikardt. Tight lower bounds for query processing on streaming and external memory data. In *Proc. International Colloquium on Automata, Languages and Programming (ICALP)*, 2005.
- [16] A. Gupta and D. Suciu. Stream processing of XPath queries with predicates. In *Proc. ACM SIGMOD International Conference on Management of Data (SIGMOD)*, pages 419–430, 2003.
- [17] D. Olteanu, T. Furche, and F. Bry. An evaluation of regular path expressions with qualifiers against XML streams. In *Proc. British National Conference on Databases (BNCD)*, pages 31–44, 2004.
- [18] P. Ramanan. Evaluating an XPath query on a streaming XML document. In *Proc. International Conference on Management of Data (COMAD)*, pages 41–52, 2005.
- [19] I. H. Sudborough. Time and tape bounded auxiliary pushdown automata. In *Proc. Mathematical Foundations of Computer Science (MFCS)*, pages 493–503. Springer Verlag, 1977.
- [20] M. Yannakakis. Algorithms for acyclic database schemes. In *Proc. International Conference on Very Large Data Bases (VLDB)*, pages 82–94, 1981.