

A Linear-Space Top-down Algorithm for Tree Inclusion Problem

Yangjun Chen¹, Yibin Chen²

Dept. Applied Computer Science, University of Winnipeg, Canada

¹y.chen@uwinnipeg.ca, ²chenyibin@gmail.com

Abstract— We consider the following tree-matching problem: Given labeled, ordered trees P and T , can P be obtained from T by deleting nodes? Deleting a node v entails removing all edges incident to v and, if v has a parent u , replacing the edges from u to v by edges from u to the children of v . The best known algorithm for this problem needs $O(|T| \cdot |\text{leaves}(P)|)$ time and $O(|\text{leaves}(P)| \cdot \min\{D_T, |\text{leaves}(T)|\} + |T| + |P|)$ space, where $\text{leaves}(T)$ (resp. $\text{leaves}(P)$) stands for the set of the leaves of T (resp. P), and D_T (resp. D_P) for the height of T (resp. P). In this paper, we present an efficient algorithm that requires $O(|T| \cdot |\text{leaves}(P)|)$ time and $O(|T| + |P|)$ space.

I. INTRODUCTION

Let T be a tree and v be a node different of root in T with parent node u . Denote by $\text{delete}(T, v)$ the tree obtained from T by removing the node v . The children of v become the children of u as illustrated in Fig. 1.

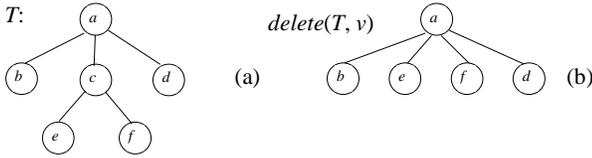


Fig. 1 The effect of removing a node from a tree

Given two ordered labeled trees P and T , called the pattern and the target, respectively. An interesting problem is: Can we obtain pattern P by deleting some nodes from target T ? That is, is there a sequence v_1, \dots, v_k of nodes such that for

$$T_0 = T \text{ and} \\ T_{i+1} = \text{delete}(T_i, v_{i+1}) \text{ for } i = 0, \dots, k - 1,$$

we have $T_k = P$? If this is the case, we say, P is included in T [9]. Such a problem is called the *tree inclusion problem*. Ordered labeled trees appear in various research fields, including programming language implementation, natural language processing, and molecular biology.

As an example [9], consider querying grammatical structures as shown in Fig. 2(a), which is the parse tree of a natural language sentence.

One might want to locate, say, those sentences that include a verb phrase containing the verb “reads” and after it a noun “book” followed by any adverb. This is exactly the sentences whose parse tree can be obtained by deleting some nodes from the tree shown in Fig. 2(a). See Fig. 2(b) for

illustration. The ordered tree inclusion problem was initially introduced by Knuth [10], where only a sufficient condition for this problem is given. The tree inclusion has been suggested as an important primitive for expressing queries on structured document databases [4, 5, 6, 12]. A structured document database is considered as a collection of parse trees that represent the structure of the stored texts and the tree inclusion is used as a means of retrieving information from them. Another application of the ordered tree matching is the video content-based retrieval. According to Rui *et al.* [14], a video can be successfully decomposed into a hierarchical tree structure, in which each node represents a scene, a group, a shot, a frame, a feature, and so on. Especially, such a tree is an ordered one since the temporal order is very important for video. In addition, the ordered tree matching can also be applied in the scene analysis, the computational biology (such as the RNA structure matching [11]), as well as in the data mining (such as the tree mining [15]).

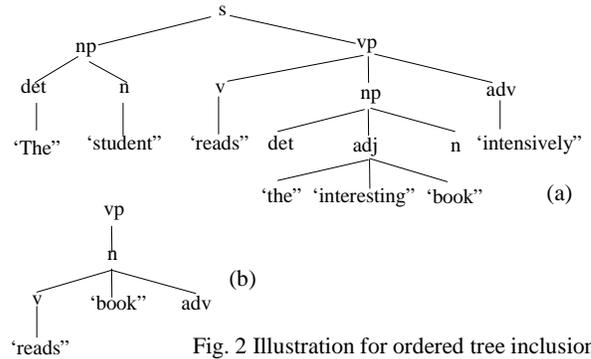


Fig. 2 Illustration for ordered tree inclusion

This problem has been the attention of much research. Kilpeläinen and Mannila [9] presented the first polynomial time algorithm using $O(|T| \cdot |P|)$ time and space. Most of the later improvements are refinements of this algorithm. In [13], Richter gave an algorithm using $O(|\alpha(P)| \cdot |T| + m(P, T) \cdot D_T)$ time, where $\alpha(P)$ is the alphabet of the labels of P , $m(P, T)$ is the size of a set called *matches*, defined as all the pairs $(v, w) \in P \times T$ such that $\text{label}(v) = \text{label}(w)$, and D_T (resp. D_P) is the depth of T (resp. P). Hence, if the number of matches is small, the time complexity of this algorithm is better than $O(|T| \cdot |P|)$. The space complexity of the algorithm is $O(|\alpha(P)| \cdot |T| + m(P, T))$. In [3], a more sophisticated algorithm was presented using $O(|T| \cdot |\text{leaves}(P)|)$ time and $O(|\text{leaves}(P)| \cdot \min\{D_T, |\text{leaves}(T)|\} + |T| + |P|)$ space. In [1], an efficient average case algorithm was discussed. Its average time complexity is $O(|T|$

+ $C(P, T) \cdot |P|$), where $C(P, T)$ represents the number of T 's nodes that have been examined during the inclusion search. However, its worst time complexity is still $O(|T| \cdot |P|)$. In [2], another bottom-up algorithm is proposed. The time complexity of the algorithm is bounded by

$$\min \begin{cases} O(|T| \cdot |\text{leaves}(P)|) \\ O(|\text{leaves}(T)| \cdot |\text{leaves}(P)| \cdot \log \log |\text{leaves}(P)| + |\text{leaves}(P)|) \\ O(|T| \cdot |P| / (\log |T|) + |T| \log |T|) \end{cases}$$

But it is claimed that the algorithm needs only $O(|T| + |P|)$ space. A careful analysis reveals that the space complexity of the algorithm is the same as that of [3]. In the algorithm, a data structure $EMB(v)$ for each v in P is used to record deep occurrences of $P[v]$ in T . It is of size $O(|\text{leaves}(T)|)$ in the worst case. $EMB(v)$ is generated recursively and works in a way similar to the concept of *shell* discussed in [3]. So the analysis of *shell* applies to $EMB(v)$'s.

In our earlier work [7], a top-down algorithm was proposed with $O(|T| + |P|)$ space requirement. However, its time complexity is not polynomial, as shown in [11].

In this paper, we revisit this issue and present a new top-down algorithm to remove any redundancy of [10]. The time complexity of the new one is bounded by $O(|T| \cdot |\text{leaves}(P)|)$. Although the time complexity of our algorithm is comparable to Chen's algorithm [3], it is more efficient than Chen's since in Chen's algorithm each node in T will be checked against, besides some internal nodes, all the leaf nodes in P . But in our algorithm, a node in T may be checked so many times only when some conditions are satisfied.

More importantly, our algorithm needs only linear space $O(|T| + |P|)$.

The tree inclusion problem on unordered trees is *NP*-complete [9] and not discussed in this paper.

II. BASIC DEFINITION

We concentrate on labeled trees that are ordered, i.e., the order between siblings is significant. Technically, it is convenient to consider a slight generalization of trees, namely forests. A forest is a finite ordered sequence of disjoint finite trees. A tree T consists of a specially designated node $root(T)$ called the root of the tree, and a forest $\langle T_1, \dots, T_k \rangle$, where $k \geq 0$. The trees T_1, \dots, T_k are the subtrees of the root of T or the immediate subtrees of tree T , and k is the outdegree of the root of T . A tree with the root t and the subtrees T_1, \dots, T_k is denoted by $\langle t; T_1, \dots, T_k \rangle$. The roots of the trees T_1, \dots, T_k are the children of t and siblings of each other. Also, we call T_1, \dots, T_k the sibling trees of each other. In addition, T_1, \dots, T_{i-1} are called the left sibling trees of T_i , and T_{i-1} the immediate left sibling tree of T_i . The root is an ancestor of all the nodes in its subtrees, and the nodes in the subtrees are descendants of the root. The set of descendants of a node v is denoted by $desc(v)$. A leaf is a node with an empty set of descendants.

Sometimes we treat a tree T as the forest $\langle T \rangle$. We may also denote the set of nodes in a forest F by $V(F)$. For example, if we speak of functions from a forest G to a forest F , we mean functions mapping the nodes of G onto the nodes of F . The size of a forest F , denoted by $|F|$, is the number of the nodes in F . The restriction of a forest F to a node v with

its descendants $desc(v)$ is called a subtree of F rooted at v , denoted by $F[v]$.

Let $F = \langle T_1, \dots, T_k \rangle$ be a forest. The preorder of a forest F is the order of the nodes visited during a preorder traversal. A preorder traversal of a forest $\langle T_1, \dots, T_k \rangle$ is as follows. Traverse the trees T_1, \dots, T_k in ascending order of the indices in preorder. To traverse a tree in preorder, first visit the root and then traverse the forest of its subtrees in preorder. The postorder is defined similarly, except that in a postorder traversal the root is visited after traversing the forest of its subtrees in postorder. We denote the preorder and postorder numbers of a node v by $pre(v)$ and $post(v)$, respectively.

Using preorder and postorder numbers, the ancestorship can be easily checked. If there is path from node u to node v , we say, u is an ancestor of v and v is a descendant of u . In this paper, by 'ancestor' ('descendant'), we mean a proper ancestor (descendant), i.e., $u \neq v$.

Lemma 1 Let v and u be nodes in a forest F . Then, v is an ancestor of u if and only if $pre(v) < pre(u)$ and $post(u) < post(v)$.

Proof. See Exercise 2.3.2-20 in [10] (page 347).

Similarly, we check the left-to-right ordering as follows.

Lemma 2 Let v and u be nodes in a forest F . v is said to be to the left of u if they are not related by the ancestor-descendant relationship and u follows v when we traverse F in preorder. Then, v is to the left of u if and only if $pre(v) < pre(u)$ and $post(v) < post(u)$.

Proof. The proof is trivial. \square

In the following, we use the postorder numbers to define an ordering of the nodes of a forest F given by $v \prec v'$ iff $post(v) < post(v')$. Also, $v \preceq v'$ iff $v \prec v'$ or $v = v'$. Furthermore, we extend this ordering with two special nodes $\perp \prec v \prec \top$. The *left relatives*, $lr(v)$, of a node $v \in V(F)$ is the set of nodes that are to the left of v and similarly the *right relatives*, $rr(v)$, are the set of nodes that are to the right of v .

The following definition is due to [9].

Definition 1 Let F and G be labeled ordered forests. We define an ordered embedding (φ, G, F) as an injective function $\varphi: V(G) \rightarrow V(F)$ such that for all nodes $v, u \in V(G)$,

- i) $label(v) = label(\varphi(v))$; (label preservation condition)
- ii) v is an ancestor of u iff $\varphi(v)$ is an ancestor of $\varphi(u)$, i.e., $pre(v) < pre(u)$ and $post(u) < post(v)$ iff $pre(\varphi(v)) < pre(\varphi(u))$ and $post(\varphi(u)) < post(\varphi(v))$; (ancestor condition)
- iii) v is to the left of u iff $\varphi(v)$ is to the left of $\varphi(u)$, i.e., $pre(v) < pre(u)$ and $post(v) < post(u)$ iff $pre(\varphi(v)) < pre(\varphi(u))$ and $post(\varphi(v)) < post(\varphi(u))$. (Sibling condition)

If there exists such an injective function from $V(G)$ to $V(F)$, we say, F includes G , F contains G , F covers G , or say, G can be embedded in F .

Fig. 3 shows an example of an ordered inclusion.

Let P and T be two labeled ordered trees. An embedding φ of P in T is said to be *root-preserving* if $\varphi(root(P)) = root(T)$. If there is a root-preserving embedding of P in T , we say that the root of T is an occurrence of P .

Fig. 3(b) also shows an example of a root preserving embedding. According to [9], restricting to root-preserving

embedding does not lose generality. In fact, what can be found by the top-down algorithm to be discussed is a root-preserving tree embedding.

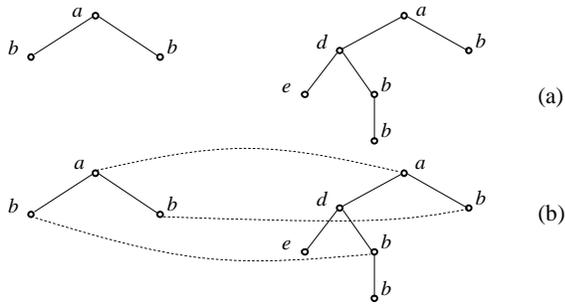


Fig. 3: (a) The tree on the left can be included in the tree on the right by deleting the nodes labelled : $d, e,$ and b ; (b) the embedding corresponding to (a).

Throughout the rest of the paper, we refer to the labeled ordered trees simply as trees.

III. ALGORITHM DESCRIPTION

Let $T = \langle t; T_1, \dots, T_k \rangle$ ($k \geq 0$) be a tree and $G = \langle P_1, \dots, P_q \rangle$ ($q \geq 0$) be a forest. We handle G as a tree $P = \langle p_v; P_1, \dots, P_q \rangle$, where p_v represent a virtual node, matching any node in T . Note that even though G contains only one single tree it is considered to be a forest. So a virtual root is added. Therefore, each node in G , except the virtual node, has a parent.

Consider a node v in $G = \langle P_1, \dots, P_q \rangle$ with children v_1, \dots, v_j . We use a pair $\langle i, v \rangle$ ($i \leq j$) to represent an ordered forest containing the first i subtrees of v : $\langle G[v_1], \dots, G[v_i] \rangle$. If v is p_v , or a node on the left-most path in P_1 , $\langle i, v \rangle$ is called a *left corner* of G . Especially, $\langle i, p_v \rangle$ is a left corner, representing the first i trees in G : P_1, \dots, P_i .

In addition, $\delta(v)$ represents a link from a node v to the left-most leaf node in $G[v]$, as illustrated in Fig. 4.

Let v' be a leaf node in G . $\delta(v')$ is defined to be a link to v' itself. So in Fig. 5, we have $\delta(v_1) = \delta(v_2) = \delta(v_3) = v_3$. We also denote by $\delta^{-1}(v')$ a set of nodes x such that for each $v \in x$ $\delta(v) = v'$. Therefore, in Fig. 5, $\delta^{-1}(v_3) = \{v_1, v_2, v_3\}$, $\delta^{-1}(v_4) = \{v_4\}$, and $\delta^{-1}(v_5) = \{v_5\}$. The out-degree of v in a tree is denoted by $d(v)$ while the height of v is denoted by $h(v)$, defined to be the number of edges on the longest downward path from v to a leaf. The height of a leaf node is set to be 0.

As with [7], we arrange two functions: *top-down*(T, G) and *bottom-up*(T', G) to check tree inclusion, where T is a tree, and T' and G are two forests. However, different from [7], each of the two functions returns a left corner $\langle i, v \rangle$ of G with the following properties:

- Let v' be the left-most leaf in $G[v]$. If $i > 0$, it shows that the first i subtrees of v in G can be embedded in T (or in T'), and for any $i' > i$, $\langle i', v \rangle$ cannot be embedded in T (or in T'), and for any v' 's ancestor $u \in \delta^{-1}(v')$ there exists no $j > 0$ such that $\langle j, u \rangle$ is able to be embedded in T (or in T').
- If $i = 0$, v is the left-most leaf in G , indicating that no left corner of G can be embedded in T (or in T').

In this sense, we say, $\langle i, v \rangle$ is the *highest* and *widest* left corner which can be embedded in T (or in T').

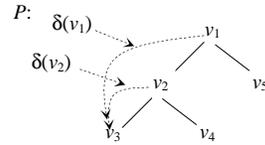


Fig. 4 A pattern tree

We notice that if $v = p_v$ and $i > 0$, it shows that P_1, \dots, P_i can be included.

In [7], both *top-down*(T, G) and *bottom-up*(T', G) return an integer i to indicate that T embeds the first i trees in G . Although our algorithm follows the arrangement of [7], the main idea is quite different. It is not necessary to refer to [7] to understand the following discussion.

If the target is a tree and the pattern is a forest, we call the function *top-down*. If both the target and the pattern are forests, we call the function *bottom-up*. But during the computation, they will be called from each other.

In *top-down*(T, G), we need to handle two cases.

Case 1: $G = \langle P_1 \rangle$; or $G = \langle P_1, \dots, P_q \rangle$ ($q > 1$), but $|T| \leq |P_1| + |P_2|$. In this case, to find the highest and widest left corner $\langle i, v \rangle$ that can be embedded in $T = \langle t; T_1, \dots, T_k \rangle$, the following checkings will be conducted:

- If t is a leaf node, we will check whether $\text{label}(t) = \text{label}(\delta(p_1))$, where p_1 is the root of P_1 . If it is the case, return $\langle 1, \text{parent of } \delta(p_1) \rangle$. Otherwise, return $\langle 0, \delta(p_1) \rangle$.
- If $|T| < |P_1|$ or $h(t) < h(p_1)$, we will make a recursive call *top-down*($T, \langle P_{11}, \dots, P_{1j} \rangle$), where $\langle P_{11}, \dots, P_{1j} \rangle$ is a forest of the subtrees of p_1 . The return value of *top-down*($T, \langle P_{11}, \dots, P_{1j} \rangle$) is used as the return value of *top-down*(T, G).
- If $|T| \geq |P_1|$ and $h(t) \geq h(p_1)$, we further distinguish between two subcases:
 - $\text{label}(t) = \text{label}(p_1)$. In this case, we will call *bottom-up*($\langle T_1, \dots, T_k \rangle, \langle P_{11}, \dots, P_{1j} \rangle$).
 - $\text{label}(t) \neq \text{label}(p_1)$. In this case, we will call *bottom-up*($\langle T_1, \dots, T_k \rangle, \langle P_1 \rangle$).

In both cases, assume that the return value of *bottom-up*() is $\langle i, v \rangle$. We need to perform a further checking:

- If $\text{label}(t) = \text{label}(v)$ and $i = d(v)$, the return value of *top-down*(T, G) is set to be $\langle 1, v \text{'s parent} \rangle$.
- Otherwise, the return value of *top-down*(T, G) is the same as $\langle i, v \rangle$.

Case 2: $G = \langle P_1, \dots, P_q \rangle$ ($q > 1$), and $|T| > |P_1| + |P_2|$. In this case, we will call *bottom-up*($\langle T_1, \dots, T_k \rangle, G$). Assume that the return value of *bottom-up*($\langle T_1, \dots, T_k \rangle, G$) is $\langle i, v \rangle$. The following checkings will be continually conducted.

- If $v \neq p_1$'s parent, check whether $\text{label}(t) = \text{label}(v)$ and $i = d(v)$. If so, the return value of *top-down*(T, G) will be set to $\langle 1, v \text{'s parent} \rangle$. Otherwise, the return value of *top-down*(T, G) is the same as $\langle i, v \rangle$.
- If $v = p_1$'s parent, the return value of *top-down*(T, G) is the same as $\langle i, v \rangle$.

The following is a formal description of the algorithm. In the process, each node t in T is associated with a data structure, referred to as $\kappa(t)$. Initially, each $\kappa(t)$ is set to ϕ . Each time a call of the form $top-down(T[t], G')$ returns a left corner $\langle i, v \rangle$, $\kappa(t)$ will be changed to $\langle i, v \rangle$, where G' is a forest made up of a set of subtrees rooted respectively at a set of consecutive child nodes (starting from a specific child to the last child) of a certain node in G . This value is mainly used in $bottom-up()$ to avoid redundancy. However, for simplicity, in the following algorithm $\kappa(t)$ is not explicitly represented.

function $top-down(T, G)$

input: $T = \langle t; T_1, \dots, T_k \rangle$, $G = \langle P_1, \dots, P_q \rangle$.

output: $\langle i, v \rangle$ specified above.

begin

```

1. if ( $q = 1$  or  $|T| \leq |P_1| + |P_2|$ )
2. then
   {let  $P_1 = \langle p_1; P_{11}, \dots, P_{1j} \rangle$ ;           (*Case 1*)
3. if  $t$  is a leaf then {
   {let  $\delta(p_1) = v$ ;                               (*Case 1 - (i)*)
4. if  $label(t) = label(v)$  then return  $\langle 1, v's\ parent \rangle$ 
   else return  $\langle 0, v \rangle$ ; }
5. if ( $|T| < |P_1|$  or  $h(t) < h(p_1)$ )
   then return  $top-down(T, \langle P_{11}, \dots, P_{1j} \rangle)$ ;   (*Case 1 - (ii)*)
6. if  $label(t) = label(p_1)$                                (*Case 1 - (iii)*)
7. then { if  $p_1$  is a leaf then {  $v := p_1's\ parent$ ;  $i := 1$ ; }
8. else {  $\langle i, v \rangle := bottom-up(\langle T_1, \dots, T_k \rangle, \langle P_{11}, \dots, P_{1j} \rangle)$ ;
9. if  $label(t) = label(v)$  and  $i = d(v)$ 
   then {  $v := v's\ parent$ ;  $i := 1$ ; }
10. }
11. else  $\langle i, v \rangle := bottom-up(\langle T_1, \dots, T_k \rangle, \langle P_1 \rangle)$ ;
   (*If  $label(t) \neq label(p_1)$ , call  $bottom-up()$ .*)
12. return  $\langle i, v \rangle$ ;
13. }
14. else
   {  $\langle i, v \rangle := bottom-up(\langle T_1, \dots, T_k \rangle, G)$ ;   (*Case 2*)
15. if  $v \neq p_1's\ parent$  then                               (*Case 2 - (iv)*)
16. if ( $label(t) = label(v)$ ) and  $i = d(v)$  then return  $\langle 1, v's\ parent \rangle$ ;
17. return  $\langle i, v \rangle$ ;                                       (*Case 2 - (v)*)
18. }
end

```

In the above algorithm, we first check whether $q = 1$ or $|T| \leq |P_1| + |P_2|$ (see line 1). If it is the case we have *Case 1* and then lines 2 - 13 are executed. In this process, all the three subcases (i), (ii), and (iii) are checked. If $q > 1$ and $|T| > |P_1| + |P_2|$, we have *Case 2* and lines 14 - 18 will be carried out, in which we first call $bottom-up(\langle T_1, \dots, T_k \rangle, G)$. Depending on its return value, (vi) or (v) is conducted.

$bottom-up(T', G)$ is designed to handle the case that both T' and G are forests made up of a set of subtrees rooted at nodes that are consecutive siblings in T and P , respectively. Let $T' = \langle T_1, \dots, T_k \rangle$ and $G = \langle P_1, \dots, P_q \rangle$. Denote by t_l the root of T_l ($l = 1, \dots, k$). Denote by p_j the root of P_j ($j = 1, \dots, q$). In $bottom-up(T', G)$, we will make a series of calls $top-down(T_l, \langle P_{j_1}, \dots, P_{j_l} \rangle)$, where $l = 1, \dots, k$, $j_1 = 1$, and $j_1 \leq j_2 \leq \dots \leq j_h \leq q$ (for some $h \leq k$), controlled as follows.

1. Two index variables l, j are used to scan T_1, \dots, T_k and P_1, \dots, P_q , respectively. (Initially, l is set to 1, and j is set to 0.) They also indicate that $\langle P_1, \dots, P_j \rangle$ has been successfully embedded in $\langle T_1, \dots, T_l \rangle$.

2. Let $\langle i_l, v_l \rangle$ be the return value of $top-down(T_l, \langle P_{j+1}, \dots, P_q \rangle)$. If $v_l = p_1's\ parent$, set j to be $j + i_l$. Otherwise, j is not changed. Set l to be $l + 1$. Go to (2).

3. The loop terminates when all $T_l's$ or all $P_j's$ are examined. If $j > 0$ when the loop terminates, $bottom-up(T', G)$ returns $\langle j, p_1's\ parent \rangle$, indicating that T' contains P_1, \dots, P_j .

Otherwise, $j = 0$, indicating that even P_1 alone cannot be embedded in any T_l ($l \in \{1, \dots, k\}$). However, in this case, we need to continue to search for a highest and widest left corner $\langle i, v \rangle$ in G , which can be embedded in T' . This is done as described below.

i) Let $\langle i_1, v_1 \rangle, \dots, \langle i_k, v_k \rangle$ be the return values of $top-down(T_1, \langle P_1, \dots, P_q \rangle), \dots, top-down(T_k, \langle P_1, \dots, P_q \rangle)$, respectively. Since $j = 0$, each $v_l \in \delta^{-1}(v')$ ($l = 1, \dots, k$), where v' is the left-most leaf in P_1 .

ii) If each $i_l = 0$, return $\langle 0, \text{left-most leaf of } P_1 \rangle$. Otherwise, there must be some $v_l's$ such that $i_l > 0$. We call such a node a *non-zero point*. Find the first non-zero point v_f with children w_1, \dots, w_s such that v_f is not a descendant of any other non-zero point. Then, we will check $\langle T_{f+1}, \dots, T_k \rangle$ against $\langle P[\lfloor w_{j+1} \rfloor], \dots, P[\lfloor w_s \rfloor] \rangle$. Let x ($0 \leq x \leq s - i_f$) be a number such that $\langle P[\lfloor w_{j+1} \rfloor], \dots, P[\lfloor w_{j+x} \rfloor] \rangle$ can be embedded in $\langle T_{f+1}, \dots, T_k \rangle$. The return value of $bottom-up(T', G)$ should be set to $\langle i_f + x, v_f \rangle$.

In the $bottom-up$ process, $\kappa(t)$ can be used to avoid redundant computation. Concretely, each time before we make a call of the form $top-down(T_l, \langle P_j, \dots, P_q \rangle)$, we will calculate a function $\kappa-checking(t_l, p_j)$ defined below to determine whether this call can be skipped over, where t_l and p_j are the roots of T_l and P_j , respectively.

function $\kappa-checking(t, p)$

input: t - a node in T ; p - a node in G .

output: ϕ or $\langle i, v \rangle$ specified above.

begin

```

1. if  $\kappa(t_l) \neq \phi$  then {
2.   let  $\kappa(t_l) = \langle i, v \rangle$ ;
3. if  $i = 0$  then return  $\phi$ ;
4. if  $i > 0$ ,  $\delta(v) = \delta(p)$ , and  $p$  is equal to  $v's\ first\ child$  or an ancestor
   of  $v's\ first\ child$ 
5. then return  $\langle i, v \rangle$ ;
6. if  $i > 0$ ,  $\delta(v) = \delta(p)$ , and  $p$  is a descendant of  $v's\ first\ child$ 
7. then return  $\langle d(p's\ parent), p's\ parent \rangle$ .
8. else return  $\phi$ .
end

```

Only when $\kappa-checking(t_l, p_j)$ returns ϕ , $top-down(T_l, \langle P_j, \dots, P_q \rangle)$ will be carried out. Otherwise, we use the value of $\kappa-checking(t_l, p_j)$ as the return value of $top-down(T_l, \langle P_j, \dots, P_q \rangle)$.

In terms of the above discussion, we arrange a new subprocedure to check a T_l against a forest $\langle P_j, \dots, P_q \rangle$, doing the same work as the $top-down$ process but with $\kappa-checking(t_l, p_j)$ being used to avoid unnecessary checkings.

function $top-down-\kappa(T, \langle P_1, \dots, P_q \rangle)$

input: T - a tree; $\langle P_1, \dots, P_q \rangle$ - a forest.

output: $\langle v, i \rangle$ specified above.

begin

1. **if** $\kappa-checking(t, p_1) = \phi$

then $\langle i, v \rangle := \text{top-down}(T, \langle P_1, \dots, P_q \rangle)$

2. **else** $\langle i, v \rangle = \kappa\text{-checking}(t, p_1)$;

3. **return** $\langle i, v \rangle$;

end

In the following algorithm, we use $\text{top-down-}\kappa(\)$, instead of $\text{top-down}(\)$, to check a tree against a forest.

function $\text{bottom-up}(T', G)$

input: $T' = \langle T_1, \dots, T_k \rangle$, $G = \langle P_1, \dots, P_q \rangle$

output: $\langle i, v \rangle$ specified above.

begin

1. $l := 1$; $j := 0$;

2. **while** $(j < q \text{ and } l \leq k)$ **do** (*main checking*)

3. $\{ \langle i_l, v_l \rangle := \text{top-down-}\kappa(T_l, \langle P_{j+1}, \dots, P_q \rangle)$

4. **if** $(v_l = p_1$'s parent and $i_l > 0)$ **then** $j := j + i_l$;

5. $l := l + 1$; $\}$

6. **if** $j > 0$ **then** return $\langle j, p_1$'s parent \rangle ;

7. **if** for all $\langle i_l, v_l \rangle$'s $i_l = 0$ **then** return $\langle 0, \text{left-most leaf in } G \rangle$

8. **else** $\{ \text{let } v_f \text{ be the first non-zero point such that it is not a descendant of any other non-zero point}$;

9. $\text{let } w_1, \dots, w_s \text{ be the children of } v_f$;

10. $l := f + 1$; $j := i_f$;

11. **while** $(j < s \text{ and } l \leq k)$ **do** (*supplement checking*)

12. $\{ \langle i_l, v_l \rangle := \text{top-down-}\kappa(T_l, \langle G[w_{j+1}], \dots, G[w_s] \rangle)$;

13. **if** $(v_l = v_f \text{ and } i_l > 0)$ **then** $j := j + i_l$;

14. $l := l + 1$; $\}$

15. **return** $\langle j, v_f \rangle$;

16. $\}$

end

In $\text{bottom-up}(T', G)$, we have two *while*-loops: one from line 2 to 5 and another from line 11 to 14. In the first *while*-loop, we check $\langle T_1, \dots, T_k \rangle$ against $\langle P_1, \dots, P_q \rangle$, referred to as the *main bottom-up checking* (or simply the *main checking*). In this checking, each T_l is checked one by one, by repeatedly calling $\text{top-down-}\kappa(T_l, \langle P_{j+1}, \dots, P_q \rangle)$ (line 3), by which $\kappa\text{-checking}(t, p_{j+1})$ is used to remove redundancy (see lines 1 - 2 in $\text{top-down-}\kappa(\)$).

In the second *while*-loop, we do a *supplement checking*. This is carried out only when the following two conditions are satisfied (see lines 6 and 7):

(1) $j = 0$, and

(2) There exists at least a non-zero point $\langle i_l, v_l \rangle$ (return value of $\text{top-down-}\kappa(T_l, \langle P_1, \dots, P_q \rangle)$) such that $i_l > 0$.

We refer to these two conditions as the *supplement checking condition*.

Let v_f be the first non-zero point such that v_f is not a descendant of any other non-zero point. Let w_1, \dots, w_s be the children of v_f . In the supplement checking, we will check $\langle T_{j+1}, \dots, T_k \rangle$ against $\langle G[w_{j+1}], \dots, G[w_s] \rangle$ (see lines 10 - 16.)

IV. CONCLUSION

In this paper, a new algorithm is proposed to improve the algorithm discussed in [7]. The main idea behind it is to let any subprocedure call return a pair to indicate a subtree (subforest) embedding while in [7], only a single integer is returned to indicate whether a whole forest (or the first several subtrees of the forest) is embedded by the corresponding target subtree. Together with a simple data structure associated with each node in the target tree to transfer the result obtained in a previous step to the next step

computation to avoid any useless effort, high performance is achieved. The time complexity of the new algorithm is bounded by $O(|T| \cdot |\text{leaves}(P)|)$ while the space requirement is bounded by $O(|T| + |P|)$, where T and P are a target and a pattern tree, respectively.

REFERENCES

- [1] L. Alonso and R. Schott. On the tree inclusion problem. In *Proceedings of Mathematical Foundations of Computer Science*, pages 211-221, 1993.
- [2] P. Bille and I.L. Gørtz. An Ordered Tree Inclusion Algorithm Based on Dynamic Tree Labeling, in *Proc.32th Intl. Colloquium on Automata, Languages and Programming*, Lecture Notes in Computer Science, vol. 3580, 2005, pp. 66-77.
- [3] W. Chen. More efficient algorithm for ordered tree inclusion. *Journal of Algorithms*, 26:370-385, 1998.
- [4] Y. Chen and Y.B. Chen. Subtree Reconstruction, Query Node Intervals and Tree Pattern Query Evaluation, *Journal of Information Science and Engineering* 28, 263-293 (2012).
- [5] Y. Chen and L. Zou, and Unordered tree matching and ordered tree matching: the evaluation of tree pattern queries, *Int. J. Information Technology, Communications and Convergence*, Vol. 1, No. 3, 2011, pp. 254-279.
- [6] Y. Chen, A New Algorithm for Twig Pattern Matching, in: *Proc. of Int. Conf. on Enterprise Information Systems (ICEIS'2007)*, IEEE, Funchal-madeira, Portugal, June 2007, pp. 44-51.
- [7] Y. Chen and Y.B. Chen, A New Tree Inclusion Algorithm, *Information Processing Letters* 98(2006) 253-262, Elsevier Science B.V.
- [8] H.L. Cheng and B.F. Wang, On Chen and Chen's new tree inclusion algorithm, *Information Processing Letters*, 2007, Vol. 103, 14-18, Elsevier Science B.V.
- [9] P. Kilpeläinen and H. Mannila. Ordered and unordered tree inclusion. *SIAM J. Comput.*, 24:340-356, 1995.
- [10] D.E. Knuth, *The Art of Computer Programming, Vol. 1 (1st edition)*, Addison-Wesley, Reading, MA, 1969.
- [11] R.B. Lyngs, M. Zuker & C.N.S. Pedersen, Internal loops in RNA secondary structure prediction, in *Proceedings of the 3rd annual international conference on computational molecular biology (RECOMB)*, 260-267 (1999).
- [12] H. Mannila and K.-J. Räihä, On Query Languages for the p-string data model, in "Information Modelling and Knowledge Bases" (H. Kangassalo, S. Ohsuga, and H. Jaakola, Eds.), pp. 469-482, IOS Press, Amsterdam, 1990.
- [13] Thorsten Richter. A new algorithm for the ordered tree inclusion problem. In *Proceedings of the 8th Annual Symposium on Combinatorial Pattern Matching (CPM)*, in *Lecture Notes of Computer Science (LNCS)*, volume 1264, pages 150-166. Springer, 1997.
- [14] Y. Rui, T.S. Huang, and S. Mehrotra, Constructing table-of-content for videos, *ACM Multimedia Systems Journal, Special Issue Multimedia Systems on Video Libraries*, 7(5):359-368, Sept 1999.
- [15] M. Zaki, Efficiently mining frequent trees in a forest. In *Proc. of KDD*, 2002.