

# Property testing of regular tree languages

Frédéric Magniez<sup>1</sup> and Michel de Rougemont<sup>2</sup>

<sup>1</sup> CNRS-LRI, UMR 8623 Université Paris-Sud, France, [magniez@lri.fr](mailto:magniez@lri.fr)

<sup>2</sup> LRI & Université Paris II, France [mdr@lri.fr](mailto:mdr@lri.fr)

**Abstract.** We consider the Edit distance with moves on the class of words and the class of ordered trees. We first exhibit a simple tester for the class of regular languages on words and generalize it to the class of ranked regular trees. We also show that the distance problem is NP-complete on ordered trees.

## 1 Introduction

Inspired by the notion of Self-testing [4, 5], Property testing has been initially defined and studied for graph properties [9]. It has been successfully extended for various classes of finite structures. Let  $\mathbf{K}$  be a class of finite structures and a distance function  $\text{dist}$ , i.e. a function between structures of  $\mathbf{K}$ . An  $\varepsilon$ -tester for a class  $\mathbf{K}_0 \subseteq \mathbf{K}$  is a randomized algorithm which takes a structure  $U_n$  of size  $n$  as input and decides if  $U_n \in \mathbf{K}_0$  or if the probability that  $U_n$  is  $\varepsilon$ -far from  $\mathbf{K}_0$  is large. A class  $\mathbf{K}_0$  is testable if for every sufficiently small  $\varepsilon$  there exists an  $\varepsilon$ -tester for  $\mathbf{K}_0$  whose time complexity is in  $O(f(\varepsilon))$ , i.e. independent of  $n$ .

For the Hamming distance, regular languages and  $\Sigma_2$ -definable graphs are testable [2, 1]. Testers have also been generalized to the infinite regular languages [7].

In this paper we initiate the study of Property testing with the Edit distance, when insertions and deletions of letters on words, of nodes and edges on trees, are the elementary operations. We specifically require an additional operation: the *move* of an entire subword or subtree in one step.

First (Section 3), we develop a new tester for regular languages on words that greatly simplifies the tester of [2] and improves its complexity by a  $\log(1/\varepsilon)$  factor.

Then (Section 4), we initiate the study of property testing on trees. The testability of regular tree languages is a well known open problem [7] for the standard Edit distance. We solve this problem when moves are allowed, by proving the testability of regular ranked tree languages.

The Word Edit distance with moves decision problem and the standard Tree Edit distance decision problem are computable in polynomial time [8, 13]. We prove in the appendix that the Tree Edit distance with moves is NP-complete. It is then interesting to point out that this apparently more complex distance yields a tester for regular languages, whereas we do not know such a tester for the classical Tree Edit distance.

Finally (Section 5), we discuss the possibility of generalizing the testability to unranked trees. As a direct application, it would imply that one can decide in constant time if a large XML document follows a DTD or is far from it.

## 2 Preliminaries

### 2.1 Property testing

Recall the notion of a (property) tester [9] on a class  $\mathbf{K}$  of finite structures for which a distance function between structures has been defined.

We say that two structures  $U_n, V_m \in \mathbf{K}$ , whose domains are respectively of size  $n$  and  $m$ , are  $\varepsilon$ -close if their distance is less than  $\varepsilon \times \max(n, m)$ . They are  $\varepsilon$ -far if they are not  $\varepsilon$ -close. In this paper, we consider this notion of closeness for words and trees since the representation of their structure is of linear size. For other classes, such as graphs, one may define the closeness relatively to the representation size (for e.g.,  $\varepsilon n^2$  for graphs) instead of the domain size.

**Definition 1.** *Let  $\varepsilon \geq 0$  be a real. An  $\varepsilon$ -tester for a class  $\mathbf{K}_0 \subseteq \mathbf{K}$  is a randomized algorithm  $A$  such that:*

- (1) *If  $U \in \mathbf{K}_0$ ,  $A$  always accepts;*
- (2) *If  $U$  is  $\varepsilon$ -far from  $\mathbf{K}_0$ , then  $\Pr[A \text{ rejects}] \geq 2/3$ .*

The *query complexity* is the number of boolean queries to the structure  $U$  of  $\mathbf{K}$ . The *time complexity* is the usual time complexity where the complexity of a query is one and the time complexity of an arithmetic operation is also one.

A class  $\mathbf{K}_0 \subseteq \mathbf{K}$  is *testable* if for every sufficiently small  $\varepsilon > 0$ , there exists an  $\varepsilon$ -tester whose time complexity depends only on  $\varepsilon$ .

### 2.2 Words

Let  $\Sigma$  be a finite alphabet of constant size and for the sake of simplicity, the reader might think that  $\Sigma = \{0, 1\}$ . We now consider the words on the alphabet  $\Sigma$ . Every word  $W$  is a finite structure  $(N, [N], l : [N] \rightarrow \Sigma)$ , where  $[N]$  denote the set  $\{1, \dots, N\}$ . The class  $\mathbf{K}$  is the set of all such structures. We will denote a subclass  $\mathbf{K}_0$  of  $\mathbf{K}$  as a subset  $L \subseteq \Sigma^*$ . In this context, a query  $i$  to some word  $W$  asks the letter  $W[i] = l(i)$ .

An *elementary operation (on words)* is a deletion or an insertion of a letter, or a move: given a subword and a position, a *move* is a transformation of the word, where the subword has been removed of its current position and inserted in the given position. Notice we omit letter replacement operations since such an operation can be simulated using one deletion and one insertion.

The standard Edit distance only considers the operations without moves, and this new distance is essential for most of the arguments.

**Definition 2.** *The distance between two words  $W$  and  $W'$  is the minimum number of elementary operations necessary to reach  $W'$  from  $W$ , noted  $\text{dist}(W, W')$ . The distance between  $W$  and a language  $L$ , noted  $\text{dist}(W, L)$ , is the minimum  $\text{dist}(W, W')$  when  $W' \in L$ .*

Let  $W$  be a word. A word  $w$  is a *subword* of  $W$  if  $w = W[i, \dots, j]$ , for some  $i, j$ .

### 2.3 Trees

Let  $T$  be an ordered  $\Sigma$ -tree, *i.e.* a tree with labels  $\sigma \in \Sigma$  on the nodes. It is *ranked* if the degree is bounded by a fixed constant, and *unranked* otherwise. We omit the term ‘ordered’, since all our trees will be ordered.

Let us first consider  $r$ -ranked trees for some fixed constant  $r$ . An  $r$ -ranked tree  $T$  is a finite structure

$$(N, [N], \text{root}, l : [N] \rightarrow \Sigma, d : [N] \rightarrow [r], s : [N] \times [r] \rightarrow [N]),$$

where  $N$  is the size of  $T$ ,  $\text{root}$  is the distinguished element representing the root of  $T$ ,  $l$  is the label function,  $d$  is the degree function which gives the degree of any node, and  $s$  is the successor partial function which associates to every node  $v$  and any position  $i \in [d(v)]$  the  $i$ -th successor of  $v$ .

The class  $\mathbf{K}$  is the set of all such structures. We will denote a subclass  $\mathbf{K}_0$  of  $\mathbf{K}$  as a subset  $L$  of all  $r$ -ranked trees. In this context, a query  $(v, i)$  to some tree  $T$  asks the label and the degree of the node  $v$  and its  $i$ -th node successor in  $T$ , if  $i \leq d(v)$ .

The classical Tree Edit distance [13] assumes basic insertions, deletions on a tree and modifications of labels (see Figure 1). A *node insertion*  $(u, \sigma)$  on an edge  $(v_1, v_2)$  replaces the edge  $(v_1, v_2)$  by the edge  $(v_1, u)$ , set  $v_2$  to be the only successor of  $u$ , and labels  $u$  by  $\sigma$ . A *node deletion* is the inverse of a node insertion. An *edge insertion*  $(v, u, \sigma, i)$  to a node  $v$  of  $T$  insert the leaf  $u$  with label  $\sigma$  between the  $(i-1)$ -th and the  $i$ -th successor of  $v$ , provided that  $d(v) < r$ . The inverse operation is an *edge deletion*.

We will also allow some moves in  $T$  (see Figure 1). A *complete subtree*  $t$  of  $T$  takes a node of  $T$  as root and is the substructure restricted to a subset of nodes such that all leaves of  $t$  are also leaves of  $T$ . A *move*  $(t, v, i)$  of a complete subtree  $t$  to a node  $v$  moves in one step  $t$  between the  $(i-1)$ -th and the  $i$ -th successor of  $v$ , provided the degree of  $v$  allows it.

An *elementary operation (on trees)* is one of the above operations.

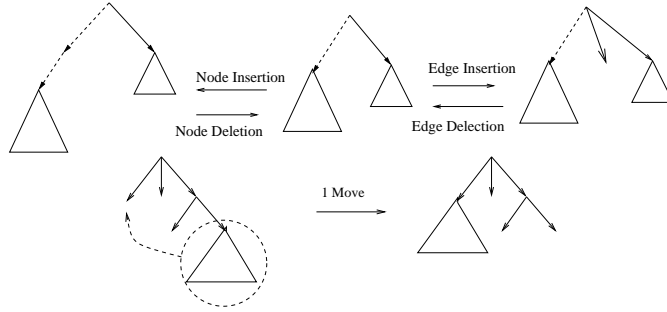
We define  $\text{dist}(T, T')$  and  $\text{dist}(T, L)$  as in Definition 2, for any trees  $T, T'$  and tree language  $L$ .

For unranked trees, the above definitions might be generalized by removing the degree condition and replacing  $[r]$  by  $[N]$ . Moreover the definition of a complete subtree is adapted so that  $t$  is a complete subtree of  $T$  if in addition it satisfies: every successors in  $t$  of every node  $v$  of  $t$  are subwords of the successors in  $T$  of  $v$ .

## 3 Testing regular languages

### 3.1 Basic definitions

Let  $\mathcal{A}$  be a deterministic automaton on words with  $m$  states,  $m \geq 2$ , which recognizes a language  $L$ . We say that  $w$  *connects* the states  $q_1$  to the state  $q_2$



**Fig. 1.** Elementary operations on trees.

when starting from  $q_1$ , the automaton  $\mathcal{A}$  reaches  $q_2$  after reading word  $w$ . If  $w$  connects  $q_1$  to  $q_2$ , we also say that  $q_1$  is *connected to*  $q_2$ . This notion will be used for random subwords  $w$  of a fixed word  $W$ .

**Proposition 1.** *Let  $q_1$  be a state connected to  $q_2$ . Then there exists a word  $w$  of size at most  $m$  that connects  $q_1$  to  $q_2$ .*

Let  $G(\mathcal{A})$  be the directed graph whose vertices are the states of  $\mathcal{A}$  and edges connects states that are connected by a word of size 1, that is a letter. We assume without loss of generality that  $G(\mathcal{A})$  is connected. Since we will only speak about strongly connected components, we omit the term ‘strongly’. A connected component  $C$  of  $G(\mathcal{A})$  is *truly connected* if there is a non empty path of  $G(\mathcal{A})$  inside  $C$ . Observe that a non truly connected component is necessarily a singleton. We will denote by  $\widehat{G}(\mathcal{A})$  the graph of the connected components of  $G(\mathcal{A})$ .

Let  $G(\mathcal{A})$  be the directed graph whose vertices are the states of  $\mathcal{A}$  and edges connects states that are connected by a word of size 1, that is a letter. We assume without loss of generality that  $G(\mathcal{A})$  is connected. Since we will only speak about strongly connected components, we omit the term ‘strongly’. A connected component  $C$  of  $G(\mathcal{A})$  is *truly connected* if there is a non empty path of  $G(\mathcal{A})$  inside  $C$ . Observe that a non truly connected component is necessarily a singleton. Let  $\widehat{G}(\mathcal{A})$  denote the graph of the connected components of  $G(\mathcal{A})$ .

**Definition 3.** *Let  $\Pi = (C_1, \dots, C_k)$  be a path of  $\widehat{G}(\mathcal{A})$ . Then  $\Pi$  is admissible if  $C_1$  (resp.  $C_k$ ) contains an initial (resp. final) state.*

**Definition 4.**

1. *Let  $C$  be a truly connected component of  $G(\mathcal{A})$ . A word  $w$  is  $C$ -simply feasible if it connects two states of  $C$ .*
2. *Let  $\Pi$  be a path of  $\widehat{G}(\mathcal{A})$ . A word  $w$  of is  $\Pi$ -feasible if it connects two states  $q_1$  and  $q_2$  along a path visiting only some of the connected components of  $\Pi$ .*

A word  $w$  is (simply)  $\Pi$ -infeasible if it is not (simply)  $\Pi$ -feasible.

A *cut* of a word  $W$  is an ordered partition of  $W$  in subwords. We will think on this partition as an ordered forest of words. Below we omit the term ‘ordered’. A cut  $F$  is  $\Pi$ -feasible if every word of  $F$  is  $\Pi$ -feasible.

### 3.2 The tester

The tester takes random subwords of finite length of  $W$  and will test feasibility for finitely many  $\Pi$ , that is at most  $2^m$  where  $m$  is the number of state of the automaton. The Robustness lemma will insure that if a word  $W$  is far, then with high probability a random subword of finite length will be infeasible.

**Tester for regular language  $(\mathcal{A}, \varepsilon, W)$ :**

If the size  $n$  of  $W$  is less than  $15m^2/\varepsilon$   
then simply evaluate  $\mathcal{A}$  on  $W$  and accept iff  $\mathcal{A}$  accepts  $W$ .  
Else do the following:

For  $i = 1, \dots, \log(5m^2/\varepsilon)$  {  
  Compute  $N_i = \Theta(\frac{2^{-i}m^3 \log(m^2/\varepsilon)}{\varepsilon})$ .  
  Choose  $N_i$  random subwords  $w_j^i$  of  $W$  of size  $2^{i+1}$ , for  $j = 1, \dots, N_i$  }  
For every admissible path  $\Pi$  of  $\widehat{G}(\mathcal{A})$  {  
  If all the  $w_j^i$  are  $\Pi$ -feasible then accept  $W$  (and stop) }  
Reject  $W$ .

**Theorem 1.** *For every real  $\varepsilon > 0$ , every automaton  $\mathcal{A}$  with  $m$  states, and every word  $W$ , the algorithm **Tester for regular language  $(\mathcal{A}, \varepsilon, W)$**  is an  $\varepsilon$ -tester for the language recognized by  $\mathcal{A}$ . Moreover, its query complexity is in  $O(m^3 \log^2(m^2/\varepsilon)/\varepsilon)$ , and its time complexity in  $O(2^m m^3 \log^2(m^2/\varepsilon)/\varepsilon)$ .*

*Proof.* We can assume w.l.o.g. that the size  $n$  of  $W$  is at least  $15m^2/\varepsilon$ , otherwise the proof of the correctness is obvious.

First, if  $W \in L$  then  $W$  is  $\Pi$ -feasible for some admissible  $\Pi$ . Therefore every subword of  $W$  is  $\Pi$ -feasible for this path  $\Pi$ . Thus the tester accepts  $W$  with probability 1.

Now suppose that  $\text{dist}(W, L) > \varepsilon n$ . Let us fixed an admissible path  $\Pi$ . Using the Robustness lemma (Lemma 1), we get that the probability to accept  $W$  for this  $\Pi$  is in  $O(2^{-m})$ . Since there is at most  $2^m$  candidates  $\Pi$ , we can conclude, using the union bound, that the overall acceptance probability is upper bounded by  $1/3$ .  $\square$

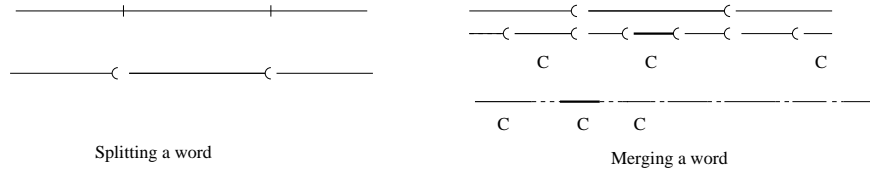
**Corollary 1.** *Regular properties of words are  $\varepsilon$ -testable.*

We now state the Robustness lemma. The notion of robustness was first defined in [12] and studied in [11]. In the rest of this section, we fix an automaton  $\mathcal{A}$  and we call  $L$  its associated language.

**Lemma 1 (Robustness).** *Let  $n \geq 15m^2/\varepsilon$ , and let  $W$  be a word of size  $n$  such that  $\text{dist}(W, L) \geq \varepsilon n$ . Then for every admissible path  $\Pi$  of  $\widehat{G}(\mathcal{A})$ , there exists an integer  $1 \leq i \leq \log(5m^2/\varepsilon)$ , such that the number of  $\Pi$ -infeasible subwords of size  $2^{i+1}$  is at least  $\frac{2^{i+1}}{90m^2 \log(5m^2/\varepsilon)} \times \varepsilon n$ .*

The sketch of the proof of the Robustness lemma takes the following steps (see Figure 2):

1. The Splitting lemma shows that if the distance between  $W$  and  $L$  is large then there are many infeasible disjoint subwords. Its proof is by contraposition:
  - (a) First, from a cut of minimal infeasible subwords, we construct a close feasible cut  $F$ .
  - (b) Then the Merging lemma which shows that if a cut  $F$  is feasible, then it is close to  $L$ .
2. The Amplifying lemma shows that if there many infeasible words, then there are many short ones.



**Fig. 2.** The correction (steps 1.a and 1.b) of a word with two infeasible subwords where  $C$  is some connected components (and  $h' = 3$  for the proof of Lemma 2).

### 3.3 Robustness of the tester

**Lemma 2 (Splitting).** *Let  $\Pi$  be an admissible path of  $\widehat{G}(\mathcal{A})$ . Let  $W$  be a word such that  $\text{dist}(W, L) > h$ . Then  $W$  has more than  $\frac{h-3m^2}{2m^2}$   $\Pi$ -infeasible disjoint subwords.*

*Proof.* The proof is by contraposition and we understand feasible as  $\Pi$ -feasible. First we construct a cut  $\mathcal{P}$  of  $W$  of size  $h'$  whose  $h' - 1$  first subwords are minimal infeasible and disjoint subwords. The last subword of  $\mathcal{P}$  is either infeasible or feasible. And in this last case, the entire word  $W$  might be feasible and  $h' = 1$ .

We visit  $W$  from the left to the right and the construction of each infeasible subword  $W[i, \dots, j]$  is done by induction on that walk.

Initially:  $h' = 0$  and  $i = j = 1$ .  
While ( $j \leq |W|$ ) {  
  While (subword  $W[i, \dots, j]$  is  $\Pi$ -feasible and  $j < |W|$ ) {increase  $j$ }  
   $h' = h' + 1$ ,  $w_{h'} = W[i, \dots, j]$ ,  
   $i = j + 1$ ,  $j = i$ .  
}

At the end of the procedure we get the desired partition  $\mathcal{P} = (w_i)_{1 \leq i \leq h'}$ .

Now we explain how to get a word  $W' \in L$ . Let  $w'_i$  be  $w_i$  without the last letter, for  $i = 1, \dots, h'$ . When  $w_{h'}$  is feasible then  $w'_{h'} = w_{h'}$ . By construction of

$w_i$ , the subwords  $w'_i$  are feasible. Let  $F$  be the cut of the  $(w'_i)_{1 \leq i \leq h'}$ . Applying Lemma 3, we get that  $\text{dist}(F, L) \leq m + 2m^2 \times h'$ . Because  $\text{dist}(W, F) \leq h'$ , then  $\text{dist}(W, L) \leq m + 2m^2 \times h'$ . But by assumption,  $h' - 1 \leq \frac{h-3m^2}{2m^2}$ , therefore  $\text{dist}(W, L) \leq h$ .  $\square$

**Lemma 3 (Merging).** *Let  $\Pi = (C_1, \dots, C_k)$  be an admissible path of  $\widehat{G}(\mathcal{A})$ . Let  $F$  be a  $\Pi$ -feasible cut of size  $h'$ . Then  $\text{dist}(F, L) \leq m + 2m^2 h'$ .*

*Proof.* First, we split each subword of  $F$  in  $C$ -feasible subwords, for some  $C \in \Pi$ . Given a  $\Pi$ -feasible subword  $w$  which connects  $p \in C_i$  to  $q \in C_j$ , we follow the automaton from  $p$  to  $q$  on  $w$ , and we delete each letter leading to a new connected component. Then the subword is cut along each deleted letter.

This technique keeps subwords that are  $C$ -feasible for some truly connected component  $C$ . Moreover, each initial subword of  $F$  splits in at most  $k$  subwords from which at most  $k$  letters are deleted, where  $k$  is less than  $m$ , where  $m$  is the number of state of the automaton. Let  $(w_i)_{1 \leq i \leq l}$  be the remaining subwords of  $F$ , where  $1 \leq l \leq m \times h'$ .

Now we explain how to move and glue the remaining subwords  $w_i$  in order to get a subword  $W' \in L$ . Let  $C'_i$  be a component of  $\Pi$  such that  $w_i$  is  $C'_i$ -feasible. Let  $p_i, q_i \in C'_i$  such that  $w_i$  connects  $p_i$  to  $q_i$ . Then, we do  $(l-1)$  moves so that the components  $C'_i$  are in the order defined by  $\Pi$ . Up to some renaming, we assume now that  $(C'_1, \dots, C'_l)$  are in the same order than  $(C_1, \dots, C_k)$ , up to some repetitions.

We glue by induction. Let  $q_0$  be an initial state of  $C_1$ , and let  $p_{l+1}$  be an accepting state of  $C_k$ . For  $i = 0$  to  $i = l$  do the following. By Proposition 1, let  $g_i$  be a word of size at most  $m$  that connects  $q_i$  to  $p_{i+1}$ . By inserting  $g_i$  between  $w_i$  and  $w_{i+1}$ , we get the word  $W' = g_0.w_1.g_1 \dots w_l.g_l$ . By construction  $W' \in L$ . In this last step, we did at most  $m \times (l+1)$  insertions.

The total number of elementary operations is less than  $(mh') + (l-1) + (m(l+1)) \leq m + 2m^2 \times h'$ , since  $l \leq mh'$  and  $m \geq 2$ .  $\square$

**Lemma 4 (Amplifying).** *Let  $\Pi$  be a path of  $\widehat{G}(\mathcal{A})$ . Let  $W$  be a word of length  $n$  with at least  $h'$   $\Pi$ -infeasible disjoint subwords. Then there exists an integer  $1 \leq i \leq \log(2n/h')$  such that the number of  $\Pi$ -infeasible subwords of size  $2^{i+1}$  is at least  $\frac{2^i(h'-4)}{6 \log(2n/h')}$ .*

*Proof.* In this proof, we understand feasible as  $\Pi$ -feasible.

Let  $w_1, \dots, w'_h$  be some infeasible disjoint subwords of  $W$ . Let  $a$  be a positive integer. For every integer  $i \geq 1$ , let  $s_i = |\{w_j : 2^{i-1} + 1 \leq |w_j| \leq 2^i\}|$ . Since we have  $|\{w_j : |w_j| > 2^a\}| \leq \frac{n}{2^a}$ , we therefore get  $\sum_{i=1}^a s_i \geq h' - \frac{n}{2^a}$ .

Take  $a = \log(2n/h')$ . Then  $\sum_{i=1}^a s_i \geq \frac{h'}{2}$ , thus there exists some  $1 \leq i \leq a$  such that  $s_i \geq \frac{h'}{2a}$ .

To lower bound the number of infeasible subwords of size  $2^{i+1}$ , we count the number of subwords of size  $2^{i+1}$  that contains a least one subword  $w_j$  whose size is in  $[2^{i-1} + 1, 2^i]$ . These subwords are also infeasible since they contain one of the infeasible subwords  $w_j$ . Note that since the subwords  $w_j$  are disjoint, each

infeasible subword of length  $2^{i+1}$  contains at most 3 of the  $w_j$  of length greater than  $2^{i-1}$ . Moreover, each infeasible subword  $w_j$  of length at most  $2^i$  is included in at least  $2^i$  subwords of length  $2^{i+1}$  (except, maybe, the two first and the two last subwords). We then get that the number of infeasible subwords of size  $2^{i+1}$  is at least  $\frac{2^i}{3} \times \frac{h'-4}{2a}$ .  $\square$

*Proof (of Lemma 1).* From the Splitting lemma with  $h = \varepsilon n$ , the word  $W$  has more than  $h' = \frac{2\varepsilon n}{5m^2}$   $\Pi$ -infeasible disjoint subwords. Now, by the Amplifying lemma, there exists an integer  $1 \leq i \leq \log(5m^2/\varepsilon)$  such that the number of  $\Pi$ -infeasible subwords of size  $2^{i+1}$  is at least  $\frac{2^i((2\varepsilon n/5m^2)-4)}{6 \log(5m^2/\varepsilon)} \geq \frac{2^{i+1}}{90m^2 \log(5m^2/\varepsilon)} \times \varepsilon n$ .  $\square$

## 4 Testing regular ranked tree languages

### 4.1 Basic definitions

A  $r$ -ranked tree automaton is a 5-tuple  $\mathcal{A} = (Q, \Sigma, \delta, (I_\sigma)_{\sigma \in \Sigma}, F)$  where  $Q$  is the set of states,  $F \subseteq Q$  is the set of accepting states,  $I_\sigma \subseteq Q$  the set of initial states for  $\sigma$ , and  $\delta : Q^{\leq r} \times \Sigma \rightarrow Q$  is the transition function.

A subtree  $t$  of  $T$  takes a node  $v_1$  of  $T$  as root and is the substructure restricted to nodes  $\{v_1, \dots, v_m\}$  where  $v_2, \dots, v_m$  are connected to  $v_1$ . The leaves of  $T$  among  $\{v_1, \dots, v_m\}$  are leaves of  $t$ , while some nodes are leaves in  $t$  but not in  $T$  and called  $*$ -nodes where the new label is  $*$ . By extension, a subtree  $t$  is a tree where some of the leaves are  $*$ -nodes.

An assignment  $\lambda$  for a subtree  $t$  determines states for its leaves such that if  $u$  is a leaf with label  $l(u)$ , then  $\lambda(u) \in I_{l(u)}$ . A run on a tree  $T$  extends  $\lambda$  on all the nodes of the subtree such that if  $u$  is a node with successors  $v_1, \dots, v_l$  where  $l \leq r$  in states  $\lambda(v_1), \dots, \lambda(v_l)$  then  $\lambda(u) = \delta(\lambda(v_1), \dots, \lambda(v_l), l(u))$ . A run *accepts* if the state of the root is in  $F$ .

Two states  $q$  and  $q'$  are *connected* if there exists a finite subtree  $t$  of size at most  $r^m$  and a run  $\lambda$  such that one leaf of  $t$  is assigned the state  $q$ , and the root of  $t$  is assigned the state  $q'$ . Let  $G(\mathcal{A})$  be the directed graph whose vertices are the states of  $\mathcal{A}$  and edges connect states that are connected by a subtree of depth 1.

We assume without loss of generality that  $G(\mathcal{A})$  is connected. We define  $\widehat{G}(\mathcal{A})$  and the notion of *truly connected* as in Section 3.1, and we omit the term ‘strongly’. We consider a set  $\Pi$  of connected components of  $G(\mathcal{A})$  and generalize the notions of  $\Pi$ -feasibility for subtrees.

**Definition 5.** *Let  $\Pi$  be a set of connected components of  $G(\mathcal{A})$ . Then  $\Pi$  is admissible if there is a pair  $(T_0, \lambda_0)$ , the witness of  $\Pi$ , such that  $\lambda_0$  is an assignment of the tree  $T_0$  which visits every connected components  $\Pi$ , and no more.*

Observe that  $T_0$  can be always chosen such that its size is at most  $r^m$ .

**Definition 6.** *Let  $\Pi$  be a set of connected components of  $G(\mathcal{A})$ . A path  $\sigma$  from a leaf to the root of a subtree  $t$  is  $\Pi$ -feasible if there exists a run which visits*

along  $\sigma$  only some connected components of  $\Pi$ . A subtree  $t$  is simply  $\Pi$ -feasible if there exists a path  $\sigma$  in  $T$  such that  $\sigma$  is  $\Pi$ -feasible. A subtree  $t$  is  $\Pi$ -feasible if there exists a run  $\lambda$  such that for all paths  $\sigma$  in  $t$ , such that  $\sigma$  is  $\Pi$ -feasible for  $\lambda$ .

A subtree  $t$  is (simply)  $\Pi$ -infeasible if it is not (simply)  $\Pi$ -feasible.

We say that two subtrees of a tree  $T$  are *disjoint* if they are node disjoint except in one node that might be both a  $*$ -node of one subtree and the root of the other subtree. A *cut* of a tree  $T$  is a partial ordered partition of  $T$  in subtrees. We will think on this partition as an ordered forest of subtrees. A *forest* of subtrees is a partial ordered set of subtrees. Below we omit the term ‘ordered’.

We naturally extend the Tree Edit distance (with moves) to forests, where the move operation can now either be applied to two subtrees of the forest or take one subtree and generate two new subtrees. Since the Tree Edit distance and the Tree Forest Edit distance are 4-equivalent (see Proposition 2), we do not distinguish them for the sake of simplicity. In other words, the Tree Forest Edit distance between trees allows for some temporarily disconnection of complete subtrees.

**Proposition 2.** *If two trees  $T, T'$  have Tree Edit distance  $h$  then their Tree Forest Edit distance is in  $[h/4, h]$ .*

A forest of subtrees is  $\Pi$ -feasible if every subtree is  $\Pi$ -feasible.

## 4.2 The tester

The tester generates random  $k$ -subtrees in the following way. A  $k$ -subtree of  $T$  from  $v$  is a subtree of  $T$  with  $v$  as a root and containing every nodes at distance at most  $k$  below  $v$ . The tester is going to select subtrees  $t_j^i$ , for  $j = 1, \dots, \Theta(\frac{mr}{\varepsilon^2})$ , of depth  $i$ , for  $i = 1, \dots, r^{2m}/\varepsilon$ , and check if they are all  $\Pi$ -feasible, for some admissible  $\Pi$ .

**Tester for regular ranked tree language  $(\mathcal{A}, \varepsilon, T)$ :**

If the size  $n$  of  $T$  is in  $O(r^{2m+1}/\varepsilon)$

then simply evaluate  $\mathcal{A}$  on  $T$  and accept iff  $\mathcal{A}$  accepts  $T$ .

Else do the following:

Compute  $N = \Theta(mr^{4m+3}/\varepsilon^2)$

For  $i = 1, \dots, 2r^{2m+1}/\varepsilon$  {

Choose  $N$  random nodes  $v_j^i$ , for  $j = 1, \dots, N$ .

Query the  $i$ -subtree  $t_j^i$  of  $T$  from  $v_j^i$ , for  $j = 1, \dots, N$  }

For every admissible set  $\Pi$  of connected component of  $G(\mathcal{A})$  {

If all the  $t_j^i$  are  $\Pi$ -feasible then accept  $T$  (and stop) }

Reject  $T$ .

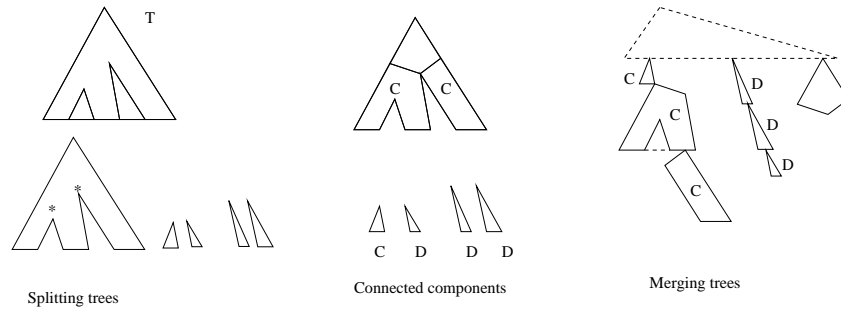
**Theorem 2.** For every real  $\varepsilon > 0$ , every  $r$ -ranked tree automaton  $\mathcal{A}$  with  $m$  states, and every  $r$ -ranked tree  $T$ , the algorithm **Tester for regular ranked tree language**  $(\mathcal{A}, \varepsilon, T)$  is an  $\varepsilon$ -tester for the language recognized by  $\mathcal{A}$ . Moreover, its query complexity is in  $O(mr^{4m+3}2^{(r^{2m+1})/\varepsilon}/\varepsilon^2)$ , and its time complexity is in  $O(2^m mr^{4m+3}2^{(r^{2m+1})/\varepsilon}/\varepsilon^2)$ .

**Corollary 2.** Regular properties of trees are  $\varepsilon$ -testable.

In the rest of this section, we fix an  $r$ -ranked automaton  $\mathcal{A}$  and we call  $L$  its associated language. The proof of Theorem 2 follow the same arguments of Theorem 1 using the Robustness lemma for trees.

**Lemma 5 (Robustness).** Let  $n = \Omega(r^{2m+1}/\varepsilon)$ , and let  $T$  be a  $r$ -ranked tree of size  $n$  such that  $\text{dist}(T, L) \geq \varepsilon n$ . Then for every admissible set  $\Pi$  of connected components of  $G(\mathcal{A})$ , there exists an integer  $1 \leq i \leq 2r^{2m+1}/\varepsilon$ , such that the number of  $\Pi$ -infeasible  $i$ -subtrees is in  $\Omega(\frac{1}{r^{4m+3}} \times \varepsilon^2 n)$ .

The structure of the proof of the Robustness lemma is the same than the one of Lemma 1 (see Figure 3).



**Fig. 3.** The correction of a tree with two infeasible subtrees where we mention  $C$  and  $D$  as some connected components (and  $h' = 3$  for the proof of Lemma 6).

### 4.3 Robustness of the tester

In this section, all the trees we consider are  $r$ -ranked trees.

**Lemma 6 (Splitting).** Let  $\Pi$  be an admissible set of connected components of the graph  $G(\mathcal{A})$ . Let  $T$  be a tree such that  $\text{dist}(T, L) > h$ . Then  $T$  has more than  $\frac{1}{3r^{m+1}}(\frac{h}{r^m} - 1) - 1$   $\Pi$ -infeasible subtrees.

*Proof.* The proof is by contraposition and we understand feasible as  $\Pi$ -feasible. First we construct a cut  $\mathcal{P}$  of  $T$  of size  $h'$  whose  $h' - 1$  last subtrees are minimal infeasible and disjoint subtrees. It might be the case that the top subtree of  $\mathcal{P}$  is  $\Pi$ -feasible. We visit  $T$  from the left to the right, and bottom-up. While visiting a node  $v$ , if the subtree below  $v$  is  $\Pi$ -infeasible, we add it in our cut

and we consider  $v$  as a  $*$ -node in the remaining part of  $T$ . At the end of the procedure we get the desired cut  $\mathcal{P} = (t_i)_{1 \leq i \leq h'}$ , ordered as  $T$  and having at most  $h'$   $*$ -node.

Now we explain how to get a tree  $T' \in L$ . Since  $t_i$  has a root of degree at most  $r$ , let  $t_i^1, \dots, t_i^r$  be the  $r$  subtrees from the root of  $t_i$  (some of them might be empty), for  $i = 1, \dots, h'$ . By construction of  $t_i$ , the subtrees  $t_i^1, \dots, t_i^r$  are  $\Pi$ -feasible. When  $t_{h'}$  is feasible then  $t_{h'}^1 = t_{h'}$  and others  $t_{h'}^j$  are empty. Let  $F$  be the forest  $(t_i^1, \dots, t_i^r)_{i=1, \dots, h'}$  of size at most  $rh'$ , in the same order than  $T$ . To get  $F$  from  $T$ , we use only  $rh'$  moves and  $h'$  edge deletions. Moreover  $F$  has at most  $h'$   $*$ -nodes. Applying Lemma 7, we get that  $\text{dist}(F, L) \leq r^m(1 + h' + 2r^{m+1}h')$  and since  $\text{dist}(T, F) \leq (r+1)h'$ , we conclude that  $\text{dist}(T, L) \leq r^m(1 + 3r^{m+1}h')$ , majoring  $(r+1) + r^m(1 + 2r^{m+1})$  by  $r^m \times 3r^{m+1}$ . But by assumption,  $h' - 1 \leq \frac{1}{3r^{m+1}}(\frac{h}{r^m} - 1) - 1$ , therefore  $\text{dist}(T, L) \leq h$ .  $\square$

**Lemma 7 (Merging).** *Let  $\Pi$  be an admissible set of connected components of  $G(\mathcal{A})$ . Let  $F$  be a  $\Pi$ -feasible forest of size  $h'_1$  with at most  $h'_2$   $*$ -nodes. Then  $\text{dist}(F, L) \leq r^m(1 + h'_2 + 2r^m h'_1)$ .*

*Proof.* First, we split each subtree  $t$  of  $F$  in simply  $C$ -feasible subtrees, for some  $C$  of  $\Pi$ . Fix such a  $t \in F$ . Let  $\lambda$  be a run of  $t$  such that all paths of  $t$  are  $\Pi$ -feasible. Fix a path  $\sigma$  of  $t$  and let  $C$  be the connected component of the root of  $t$ . We follow  $\sigma$  top-down until we leave  $C$  after a node  $v$ . Then we cut  $t$  just before leaving  $C$ , that is between  $v$  and its successors using  $r$  edge deletions and  $r$  moves. This leads to one simply  $C$ -feasible subtree from the root of  $t$  where the label of  $v$  is now  $*$ , and  $r$  new  $\Pi$ -feasible subtrees from the successors of  $v$ . We iterate the argument for the last  $r$  subtrees using the restrictions of the same run  $\lambda$ , so that the next paths of the last  $r$  subtrees will start with the next connected component. At the end of the process, at most  $1 + r + \dots + r^{m-1} \leq r^m$   $C$ -feasible subtrees have been generated from  $t$  using  $r^m$  edge deletions and  $r^m$  moves.

We only consider subtrees that are simply  $C$ -feasible for some truly connected component  $C$  of  $\Pi$ , and delete the other ones, necessarily of size 1, using at most  $r^m \times h'_1$  node deletions. Let  $(t_i)_{1 \leq i \leq k}$  be the remaining subtrees of  $F$ , where  $1 \leq k \leq r^m \times h'_1$ .

We now explain how to move and glue the remaining subtrees  $w_i$  in order to get a tree  $T' \in L$ . Let  $C_i$  be a connected component of  $\Pi$  such that  $t_i$  is simply  $C_i$ -feasible. We first move and glue linearly each subtrees  $t_i$  with the same  $C_i$ . At each  $*$ -node, a tree of size  $r^m$  is also inserted so that the resulting subtree is simply  $C_i$ -feasible and without any  $*$ -nodes. Then the remaining subtrees are connected to  $T_0$  in order to get a tree  $T' \in L$ . We have done  $(k-1)$  moves and  $r^m \times (k+1) + r^m \times h'_2$  insertions.

The total number of elementary operations is less than

$$\begin{aligned} & (2r^m \times h'_1) + (r^m \times h'_1) + (r^m \times h'_1 + r^m \times (r^m \times h'_1 + 1) + r^m \times h'_2) \\ & \leq r^m + 2r^{2m} \times h'_1 + r^m \times h'_2. \end{aligned}$$

$\square$

**Lemma 8 (Amplifying).** *Let  $\Pi$  be an admissible set of connected components of  $G(\mathcal{A})$ . Let  $T$  be a tree of size  $n$  with at least  $h'$   $\Pi$ -infeasible disjoint subtrees. Then there exists an integer  $1 \leq i \leq 2n/h'$  such that the number of  $\Pi$ -infeasible  $i$ -subtrees is at least  $\frac{1}{r} \times \frac{h'}{4n} \times h'$ .*

*Proof.* In this proof, we understand feasible as  $\Pi$ -feasible and we follow the structure of the proof of Lemma 4.

Let  $t_1, \dots, t_{h'}$  be some infeasible disjoint subtrees of  $T$ . Let  $a$  be a positive integer. For every integer  $i \geq 1$ , let  $s_i = |\{t_j : \text{depth}(t_j) = i\}|$ . Since the root of a subtree may be shared with the leaf of another subtree as a  $*$ -node, we have  $|\{t_j : \text{depth}(t_j) > a\}| \leq \frac{n}{a+1}$ , and therefore  $\sum_{i=1}^a s_i \geq h' - \frac{n}{a}$ .

Take  $a = \frac{2n}{h'}$ . Then  $\sum_{i=1}^a s_i \geq \frac{h'}{2}$ , thus there exists some  $1 \leq i \leq a$  such that  $s_i \geq \frac{h'}{2a}$ .

To lower bound the number of infeasible  $i$ -subtrees, we count the number of  $i$ -subtrees that contains a least one subtree  $t_j$  of depth  $i$ . These subtrees are also infeasible since they contain one of the infeasible subtrees  $t_j$ . Note that since the subtrees  $t_j$  are disjoint, each infeasible  $i$ -subtrees contains at most  $r$  of the  $t_j$  of depth  $i$ . Moreover, each infeasible subtree  $t_j$  of depth  $i$  is included in at least one infeasible  $i$ -subtree. We then get that the number of infeasible  $i$ -subtrees is at least  $\frac{1}{r} \times \frac{h'}{2a}$ .  $\square$

## 5 Extension to unranked trees

An *unranked tree automaton* generalizes the transition function to  $\delta : Q \times \Sigma \rightarrow 2^{Q^*}$  such that  $\delta(q, a)$  is a regular language on  $Q$ . A run  $\lambda$  is generalized such that if  $u$  is a node with successors  $v_1, \dots, v_l$  in states  $\lambda(v_1), \dots, \lambda(v_l)$  and there is a  $q$  such that  $\lambda(v_1), \dots, \lambda(v_l) \in \delta(q, l(u))$ , then  $\lambda(u) = q$ .

We consider two approaches to generalize **Tester for regular ranked tree language** to unranked regular trees. In a direct approach we are able to prove a Splitting lemma and a Merging lemma for any unranked tree automaton. The remaining main obstacle is the existence of an efficient random generator of subtrees for a corresponding Amplifying lemma.

Another possible approach consists to encode unranked trees  $T$  by binary trees  $e(T)$  using the MSV encoding (see Section A), to construct a binary automaton that accepts the encoded unranked trees accepted by the unranked automaton, and to apply the tester on binary trees. In case of XML files, assume they are given by their DOM (Document Object Model) structures. We can efficiently generate any random  $k$ -subtrees on the encoded tree from the DOM and simulate efficiently **Tester for regular ranked tree language** on the encoded tree. There is a remaining obstacle consisting in lower bounding the distance of two encoded trees  $\text{dist}(e(T), e(T'))$  by  $\text{dist}(T, T')$ . Even if it is clear that  $\text{dist}(e(T), e(T')) \leq 2\text{dist}(T, T')$ , the opposite inequality is rather technical.

## References

1. N. Alon, E. Fischer, M. Krivelevich, and M. Szegedy. Efficient testing of large graphs. *Combinatorica*, 20:451–476, 2000.
2. N. Alon, M. Krivelevich, I. Newman, and M. Szegedy. Regular languages are testable with a constant number of queries. *SIAM Journal on Computing*, 30(6), 2000.
3. A. Apostolico and Z. Galil. *Pattern Matching Algorithms*, chapter 14: Approximate Tree Pattern Matching. Oxford University Press, 1997.
4. M. Blum and S. Kannan. Designing programs that check their work. *Journal of the ACM*, 42(1):269–291, 1995.
5. M. Blum, M. Luby, and R. Rubinfeld. Self-testing/correcting with applications to numerical problems. *Journal of Computer and System Sciences*, 47(3):549–595, 1993.
6. A. Bruggemann-Klein. and D. Wood. One-unambiguous regular languages. *Information and Computation*, 142:182–206, 1998.
7. H. Chockler and O. Kupferman.  $\omega$ -regular languages are testable with a constant number of queries. In *Proceedings of the 6th Workshop on Randomization and Approximation Techniques in Computer Science*, pages 26–38, 2002. LNCS volume 2483.
8. G. Cormode. *Sequence Distance Embeddings*. PhD thesis, University of Warwick, 2003.
9. O. Goldreich, S. Goldwasser, and D. Ron. Property testing and its connection to learning and approximation. *Journal of the ACM*, 45(4):653–750, 1998.
10. T. Milo, D. Suci, and V. Vianu. Typechecking for xml transformers. *Proceedings of the ACM Symposium on Principles of Database Systems*, pages 11–22, 2000.
11. R. Rubinfeld. On the robustness of functional equations. *SIAM Journal on Computing*, 28(6):1972–1997, 1999.
12. R. Rubinfeld and M. Sudan. Robust characterizations of polynomials with applications to program testing. *SIAM Journal on Computing*, 25(2):23–32, 1996.
13. K. Tai. The tree-to-tree correction problem. *Journal of the ACM*, 26:422–433, 1979.

## A Encoding of unranked trees

Every unranked tree can be coded as a binary tree but many encodings are possible. We consider the MSV-encoding [10], which follows the following schema where  $(T)$  denotes an unranked tree and  $(F)$  an unranked forest:

$$\begin{aligned} T &:= a(F), & a \in \Sigma & & F &:= T, F \\ T &:= a(), & a \in \Sigma & & F &:= T \end{aligned}$$

The encoding into a binary tree over  $\Sigma' = \Sigma \cup \{-, |\}$  where  $-$ ,  $|$  are new symbols uses two functions  $encode$ , the encoding of a tree and  $encode_f$ , the encoding of a forest. Its recursive definition is given by:

$$\begin{aligned} encode[a(F)] &:= a(encode_f[(F)], |) & encode_f[T, F] &:= -(encode[T], encode_f[F]) \\ encode[a()] &:= a(|, |) & encode_f[T] &:= encode[T] \end{aligned}$$

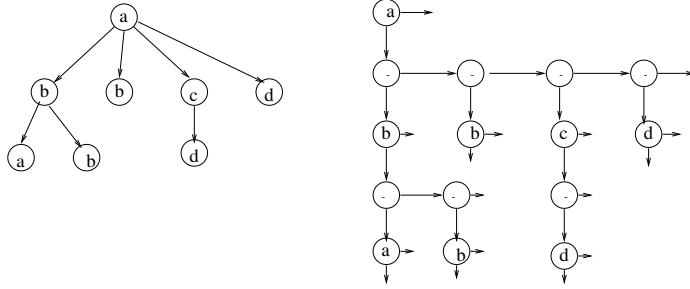


Fig. 4. MSV encoding of a tree.

It is a classical observation that a language  $L$  of unranked trees is regular iff the language  $L'$  of its MLV-encoding is regular. In the case of XML,  $Q = \Sigma$  and we associate with each  $q$  a regular expression  $r_q$  and a node labelled  $q$  is *parsed correctly* if the labels of its successors belong to the regular expression  $r_q$ . The DTD gives the sequence of rules  $q : r(q)$  but it is assumed to be 1-unambiguous [6]. In general, we associate with a state  $q$  and symbol  $a$  a regular expression  $r_{a,q}$ .

## B The Tree Edit distance with moves problem

The distance problem takes two trees  $T_1, T_2$  and an integer  $p$  as input and decides if  $\text{dist}(T_1, T_2) \leq p$ .

**Proposition 3.** *The distance problem on ordered trees with the Edit distance with moves is NP-complete on ranked and unranked trees.*

*Proof.* We reduce an NP-complete problem 3SM (3-Set Matching or Exact Cover by 3-Sets) to the distance problem on unranked trees. We use unranked trees similar to the ones mentioned in [3] where it is shown that the distance problem on unordered trees is NP-complete. Consider a set  $U = \{u_1, u_2, \dots, u_{3k}\}$  and  $n$  sets of 3 elements of  $U$ , i.e.  $S_1 = \{u_{1_1}, u_{1_2}, u_{1_3}\}, \dots, S_i = \{u_{i_1}, u_{i_2}, u_{i_3}\}, \dots, S_n = \{u_{n_1}, u_{n_2}, u_{n_3}\}$ . We have to decide if there are  $k$  subsets  $S_{i_1}, \dots, S_{i_k}$  which partition  $U$ , i.e. whose intersection is empty and their union is  $U$ .

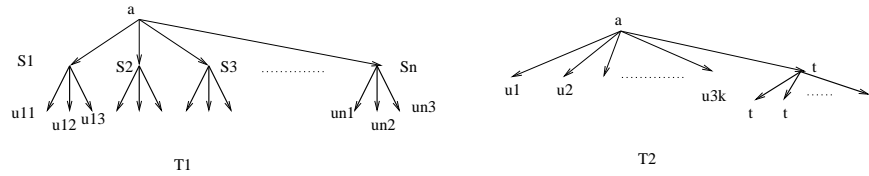
Consider the two unranked trees below,  $T_1$  and  $T_2$  on an alphabet  $\Sigma = \{u_1, u_2, \dots, u_{3k}, S_1, \dots, S_n, a, t\}$ . The tree  $T_1$  is built from the  $n$  sets  $S_1 = \{u_{1_1}, u_{1_2}, u_{1_3}\}, \dots, S_i = \{u_{i_1}, u_{i_2}, u_{i_3}\}, \dots, S_n = \{u_{n_1}, u_{n_2}, u_{n_3}\}$  as in Figure 5 and the distance between the node  $S_i$  and the three leaves is  $n + 5k$ . The tree  $T_2$  has  $u_1, \dots, u_{3k}$  as successors and a node labelled  $t$  with  $3(n - k)$  leaves at distance  $n + 5k$  labelled  $t$ .

If we can partition  $U$  with  $k$  such sets, we can use  $k$  delete of edges  $(a, S_i)$  followed by  $k$  moves to place the  $u_{i_1}, u_{i_2}, u_{i_3}$  to the left. We reorder all the  $u_{i_j}$  and use up to  $3k$  moves to obtain the left part of  $T_2$ . We then use  $n - k$  merges on the other sets to obtain the right part of  $T_2$  and  $3(n - k) + 1$  modifications of

labels after we change all the labels to  $t$ . In the worst-case we make  $2k + 3k + (n - k) + 3(n - k) + 1 = 4n + 2k + 1 = p$  operations.

If the instance of 3SM is positive, the distance between  $T_1$  and  $T_2$  is less than  $p$ . If the instance of 3SM is negative, let us show that the distance is greater than  $p$ . The transformation of the second part of the tree requires  $3(n - k) + 1$  operations. If we select any subset of size  $k$ , there will be at least one duplicate  $u_i$  which requires  $n + 5k$  operations to remove and in the best case, we apply  $k + n + 5k$  delete. The distance is at least  $4n + 3k + 1$ . This proves that the distance problem on unranked trees with an infinite alphabet is NP-complete.

To reduce the problem to a finite alphabet by replacing nodes labelled by a letter of  $\Sigma = \{u_1, u_2, \dots, u_{3k}\}$  by a binary tree of size  $\log(3k)$  and suppressing the labels  $S_1, \dots, S_n$ . Construct larger trees  $T'_1$  and  $T'_2$  from  $T_1$  and  $T_2$  by replacing the nodes labelled  $u_{i_1}, u_{i_2}, u_{i_3}$  by such finite trees. We keep the same construction but instead of  $3(n - k) + 1$  modifications of labels, we first need to apply  $3(n - k) \log(3k)$  deletions of edges to remove the finite binary trees, and then apply the  $3(n - k) + 1$  modifications of labels. We just take  $p' = 2k + 3k + (n - k) + (3(n - k) + 1) \log(3k) = (3(n - k) + 1) \log(3k) + n + 4k$ . This concludes that the distance problem on unranked trees is NP-complete. If we code the trees  $T'_1$  and  $T'_2$  by binary trees  $T''_1$  and  $T''_2$  using the MSV-encoding, we notice that if  $\text{dist}(T'_1, T'_2) = k$  then  $\text{dist}(T''_1, T''_2) = 2k$ : we reduce the distance problem on unranked trees to the distance problem on ranked trees.  $\square$



**Fig. 5.** Trees  $T_1$  and  $T_2$ .

The argument can be also used to show that the distance problem is not approximable as we can maintain an arbitrarily large gap between positive and negative instances.