Revisiting the tree edit distance and its backtracing: A tutorial

Benjamin Paaßen

Supplementary material for the ICML 2018 paper: Tree Edit Distance Learning via Adaptive Symbol Embeddings

Abstract

Almost 30 years ago, Zhang and Shasha (1989) published a seminal paper describing an efficient dynamic programming algorithm computing the tree edit distance, that is, the minimum number of node deletions, insertions, and replacements that are necessary to transform one tree into another. Since then, the tree edit distance has been widely applied, for example in biology and intelligent tutoring systems. However, the original paper of Zhang and Shasha can be challenging to read for newcomers and it does not describe how to efficiently infer the optimal edit script.

In this contribution, we provide a comprehensive tutorial to the tree edit distance algorithm of Zhang and Shasha. We further prove metric properties of the tree edit distance, and describe efficient algorithms to infer the cheapest edit script, as well as a summary of all cheapest edit scripts between two trees.

ipts between two trees.

A reference implementation of the algorithms presented in this work can be found at https://openresearch.cit-ec.de

1 Introduction

The tree edit distance (TED, Zhang and Shasha 1989) between two trees \bar{x} and \bar{y} is defined as the minimum number of nodes that need to be replaced, deleted, or inserted in \bar{x} to obtain \bar{y} . This makes the TED an intuitive notion of distance, which has been applied in a host of different application areas (Pawlik and Augsten 2011), for example to compare RNA secondary structures and phylogenetic trees in biology (Akutsu 2010; S. Henikoff and J. G. Henikoff 1992; McKenna et al. 2010; Smith and Waterman 1981), or to recommend edits to students in intelligent tutoring systems (Choudhury, Yin, and Fox 2016; Freeman, Watson, and Denny 2016; Nguyen et al. 2014; Paaßen et al. 2018; Rivers and Koedinger 2015). As such, the TED has certainly stood the test of time and is still of great interest to a broad community. Unfortunately, though, a detailed tutorial on the TED seems to lack, such that users tend to treat it as a black box. This is unfortunate as the TED lends itself for straightforward adjustments to the application domain at hand, and this potential remains under-utilized.

This contribution is an attempt to provide a comprehensive tutorial to the TED, enabling users to implement it themselves, adjust it to their needs, and compute not only the distance as such but also the optimal edits which transform \bar{x} into \bar{y} . Note that we focus here on the original version of the TED with a time complexity of $\mathcal{O}(m^2 \cdot n^2)$ and a space complexity of $\mathcal{O}(m \cdot n)$, where m and n are the number of nodes in \bar{x} and \bar{y} respectively (Zhang and Shasha 1989). Recent innovations have improved the worst-case time complexity to cubic time (Pawlik and Augsten 2011; Pawlik and Augsten 2016), but require deeper knowledge regarding tree decompositions. Furthermore, the practical runtime complexity of the original TED algorithm is still competitive for balanced trees, such that we regard it as a good choice in many practical scenarios (Pawlik and Augsten 2011).

Our tutorial roughly follows the structure of the original paper of Zhang and Shasha (1989), that is, we start by first defining trees (section 2.1) and edit scripts on trees (section 2.2), which are the basis for the TED. To make the TED more flexible, we introduce generalized cost functions on edits (section 2.3), which are a good interface to adjust the TED for custom applications. We conclude the theory section by introducing mappings between subtrees (section 2.4), which constitute the interface for an efficient treatment of the TED.

These concepts form the basis for our key theorems, namely that the cheapest mapping between two trees can be decomposed via recurrence equations, which in turn form the basis for Zhang and Shasha's dynamic programming algorithm for the TED (section 3). Finally, we conclude this tutorial with a section on the backtracing for the TED, meaning that we describe how to efficiently compute the cheapest edit script transforming one tree into another (section 4).

e cheapest edit script transforming one tree into another (section 4).

A reference implementation of the algorithms presented in this work can be found at https://openresearch.cit-e

2 Theory and Definitions

We begin our description of the tree edit distance (TED) by defining trees, forests, and tree edits, which provides the basis for our first definition of the TED. We will then revise this definition by permitting customized costs for each edit, which yields a generalized version of the TED. Finally, we will introduce the concept of a tree mapping, which will form the basis for the dynamic programming algorithm.

2.1 Trees

Definition 1 (Alphabet, Tree, Label, Children, Leaf, Subtree, Parent, Ancestor, Forest). Let \mathcal{X} be some arbitrary set which we call an *alphabet*.

We define a tree \bar{x} over the alphabet \mathcal{X} recursively as $x(\bar{x}_1,\ldots,\bar{x}_R)$, where $x\in\mathcal{X}$, and $\bar{x}_1,\ldots,\bar{x}_R$ is a (possibly empty) list of trees over \mathcal{X} . We denote the set of all trees over \mathcal{X} as $\mathcal{T}(\mathcal{X})$.

We call x the *label* of \bar{x} , also denoted as $\nu(\bar{x})$, and we call $\bar{x}_1, \ldots, \bar{x}_R$ the *children* of \bar{x} , also denoted as $\bar{\varrho}(\bar{x})$. If a tree has no children (i.e. R=0), we call it a *leaf*. In terms of notation, we will generally omit the brackets for leaves, i.e. x is a notational shorthand for x().

We define a *subtree* of \bar{x} as either \bar{x} itself, or as a subtree of a child of \bar{x} , i.e. as the transitive closure over the child relation, including the tree itself. Conversely, we call \bar{x} the *parent* of \bar{y} if \bar{y} is a child of \bar{x} , and we call \bar{x} an *ancestor* of \bar{y} if \bar{x} is either the parent of \bar{y} or an ancestor of the parent of \bar{x} , i.e. the transitive closure over the parent relation. We call the multi-set of labels for all subtrees of a tree the *nodes* of the tree.

We call a list of trees $\bar{x}_1, \dots, \bar{x}_R$ from $\mathcal{T}(\mathcal{X})$ a forest over \mathcal{X} , and we denote the set of all possible forests over \mathcal{X} as $\mathcal{T}(\mathcal{X})^*$. We denote the empty forest as ϵ .

As an example, consider the alphabet $\mathcal{X} = \{a,b\}$. Some example trees over \mathcal{X} are a,b,a(a),a(b),b(a,b), and a(b(a,b),b).

An example forest over this alphabet is a, b, b(a, b). Note that each tree is also a forest. This is important as many of our proofs in this paper will be concerned with forests, and these proofs apply to trees as well.

Now, consider the example tree $\bar{x} = a(b(c,d),e)$ from Figure 1 (left). a is the label of \bar{x} , and b(c,d) and e are the children of \bar{x} . Conversely, \bar{x} is the parent of b(c,d) and e. The leaves of \bar{x} are c, d, and e. The subtrees of \bar{x} are \bar{x} , b(c,d), c, d, and e, and the nodes of \bar{x} are a, b, c, d, and e.

2.2 Tree Edits

Next, we shall consider *edits* on trees, that is, functions which *change* trees (or forests). In particular, we define:

Definition 2 (Tree Edit, Edit Script). A tree edit over the alphabet \mathcal{X} is a function δ which maps a forest over \mathcal{X} to another forest over \mathcal{X} , that is, a tree edit δ over \mathcal{X} is any kind of function δ : $\mathcal{T}(\mathcal{X})^* \to \mathcal{T}(\mathcal{X})^*$.

In particular, we define a *deletion* as the following function del.

$$del(\epsilon) := \epsilon$$
$$del(\bar{x}_1, \dots, \bar{x}_R) := \bar{\varrho}(\bar{x}_1), \bar{x}_2, \dots, \bar{x}_R$$

We define a replacement with node $y \in \mathcal{X}$ as the following function rep_y .

$$\operatorname{rep}_y(\epsilon) := \epsilon$$

$$\operatorname{rep}_y(\bar{x}_1, \dots, \bar{x}_R) := y(\bar{\varrho}(\bar{x}_1)), \bar{x}_2, \dots, \bar{x}_R$$

And we define an insertion of node $y \in \mathcal{X}$ as parent of the trees l to r-1 as the following function $\inf_{y,l,r}$.

$$\operatorname{ins}_{y,l,r}(\bar{x}_1,\ldots,\bar{x}_R) := \begin{cases} \bar{x}_1,\ldots,\bar{x}_R & \text{if } r > R+1, l > r, \text{ or } l < 1\\ \bar{x}_1,\ldots,\bar{x}_{l-1},y,\bar{x}_l,\ldots,\bar{x}_R & \text{if } l = r \leq R+1\\ \bar{x}_1,\ldots,\bar{x}_{l-1},y(\bar{x}_l,\ldots,\bar{x}_{r-1}),\bar{x}_r,\ldots,\bar{x}_R & \text{if } 1 \leq l < r \leq R+1 \end{cases}$$

We define an *edit script* $\bar{\delta}$ as a list of tree edits $\delta_1, \ldots, \delta_T$. We define the application of an edit script $\bar{\delta} = \delta_1, \ldots, \delta_T$ to a tree \bar{x} as the composition of all edits, that is: $\bar{\delta}(\bar{x}) := \delta_1 \circ \ldots \circ \delta_T(\bar{x})$, where \circ denotes the contravariant composition operator, i.e. $f \circ g(x) := g(f(x))$.

Let Δ be a set of tree edits. We denote the set of all possible edit scripts using edits from Δ as Δ^* . We denote the empty script as ϵ .

As an example, consider the alphabet $\mathcal{X} = \{a,b\}$ and the edit rep_b , which replaces the first node in a forest with a b. If we apply this edit to the example tree a(b,a), we obtain $\operatorname{rep}_b(a(b,a)) = b(b,a)$. Now, consider the edit script $\bar{\delta} := \operatorname{del}, \operatorname{rep}_a, \operatorname{ins}_{a,1,3}$, which yields the following result for the example tree a(b,a).

$$\begin{split} \bar{\delta}(\mathbf{a}(\mathbf{b},\mathbf{a})) &= \mathrm{del} \circ \mathrm{rep}_{\mathbf{a}} \circ \mathrm{ins}_{\mathbf{a},1,3}(\mathbf{a}(\mathbf{b},\mathbf{a})) \\ &= \mathrm{rep}_{\mathbf{a}} \circ \mathrm{ins}_{\mathbf{a},1,3}(\mathbf{b},\mathbf{a}) \\ &= \mathrm{ins}_{\mathbf{a},1,3}(\mathbf{a},\mathbf{a}) \\ &= \mathbf{a}(\mathbf{a},\mathbf{a}) \end{split}$$

Note that tree edits are defined over forests, not only over trees. This is necessary because, as in our example above, deletions may change trees into forests and need to be followed up with insertions to obtain a tree again.

Based on edit scripts, we can define the TED.

Definition 3 (Edit Distance). Let \mathcal{X} be an alphabet and Δ be a set of tree edits over \mathcal{X} . Then, the TED according to Δ is defined as the function

$$d_{\Delta}: \mathcal{T}(\mathcal{X}) \times \mathcal{T}(\mathcal{X}) \to \mathbb{N} \tag{1}$$

$$d_{\Delta}(\bar{x}, \bar{y}) = \min_{\bar{\delta} \in \Delta^*} \left\{ |\bar{\delta}| \left| \bar{\delta}(\bar{x}) = \bar{y} \right. \right\}$$
 (2)

In other words, we define the TED between two trees \bar{x} and \bar{y} as the minimum number of edits we need to transform \bar{x} to \bar{y} .

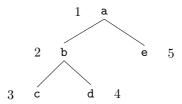
Our definition of tree edit is very broad and includes many edits which are not meaningful in most tasks. Therefore, the standard TED of Zhang and Shasha (1989) is restricted to the three kinds of special edits listed above, namely deletions, which remove a single node from a forest, insertions, which insert a single node into a forest, and replacements, which replace a single node in a forest with another node. Up to now, we have only defined versions of these edits which apply to the first node in a forest. We now go on to define variants which can be applied to any node in a given forest. To this end, we need a way to uniquely identify and target single nodes in a forest. We address this problem via the concept of a pre-order. The pre-order just lists all subtrees of a forest recursively, starting with the first tree in its forest, followed by the pre-order of its children and the pre-order of the remaining trees. More precisely, we define the pre-order as follows.

Definition 4 (Pre-Order). Let $\bar{x}_1, \ldots, \bar{x}_R$ be a forest over some alphabet \mathcal{X} . Then, we define the *pre-order* of $\bar{x}_1, \ldots, \bar{x}_R$ as the list $\pi(\bar{x}_1, \ldots, \bar{x}_R)$ which enumerates all subtrees of $\bar{x}_1, \ldots, \bar{x}_R$ as follows.

$$\pi(\epsilon) := \epsilon \tag{3}$$

$$\pi(\bar{x}_1, \dots, \bar{x}_R) := \bar{x}_1 \oplus \pi(\bar{\rho}(\bar{x}_1)) \oplus \pi(\bar{x}_2, \dots, \bar{x}_R) \tag{4}$$

where \oplus denotes list concatenation.



i	$ar{x}_i$	x_i	$p_{ar{x}}(i)$	$r_{ar{x}}(i)$
1	a(b(c,d),e)	a	0	1
2	b(c,d)	b	1	1
3	С	С	2	1
4	d	d	2	2
5	е	е	1	2

Figure 1: Left: The tree $\bar{x} = a(b(c,d), e)$ with pre-order indices drawn next to each node. Right: A table listing the subtrees \bar{x}_i , the nodes x_i , the parents $p_{\bar{x}}(i)$, and the child indices $r_{\bar{x}}(i)$ for all pre-order indices $i \in \{1, ..., 5\}$ for the tree $\bar{x} = a(b(c,d),e)$.

As a shorthand, we denote the *i*th subtree $\pi(X)_i$ as \bar{x}_i . We define x_i as the label of \bar{x}_i , i.e. $x_i := \nu(\bar{x}_i)$.

We define the *size* of the forest X as the length of the pre-order, that is $|X| := |\pi(X)|$.

Further, we define $p_{\bar{x}}(i)$ as the pre-order index of the parent of \bar{x}_i , that is, $\bar{x}_{p_{\bar{x}}(i)}$ is the parent of \bar{x}_i . If there is no parent, we define $p_{\bar{x}}(i) := 0$.

Finally, we define $r_{\bar{x}}(i)$ as the child index of \bar{x}_i , that is, $\bar{x}_i = \bar{\varrho}(\bar{x}_{p_{\bar{x}}(i)})_{r_{\bar{x}}(i)}$. If $p_{\bar{x}}(i) = 0$, we define $r_{\bar{x}}(i)$ as the index of \bar{x}_i in the forest, that is, the $r_{\bar{x}}(i)$ th tree in X is \bar{x}_i .

Consider the example of the tree $\bar{x} = a(b(c,d),e)$ from Figure 1 (left). Here, the pre-order is $\pi(\bar{x}) = a(b(c,d),e)$, b(c,d), c,d, e. Figure 1 (right) lists for all i the subtrees \bar{x}_i , the nodes x_i , the parents $p_{\bar{x}}(i)$, and the child indices $r_{\bar{x}}(i)$.

Based on the pre-order, we can specify replacements, deletions, and insertions as follows:

Definition 5 (Replacements, Deletions, Insertions). Let \mathcal{X} be some alphabet. We define a deletion of the *i*th node as the following function del_i .

$$\begin{aligned}
\det_{i}(\bar{x}_{1}, \dots, \bar{x}_{R}) &:= \begin{cases}
\bar{x}_{1}, \dots, \bar{x}_{R} & \text{if } i < 1 \\
\det(\bar{x}_{1}, \dots, \bar{x}_{R}) & \text{if } i = 1 \\
\nu(\bar{x}_{1}) \left(\det_{i-1}(\bar{\varrho}(\bar{x}_{1}))\right), \bar{x}_{2}, \dots, \bar{x}_{R} & \text{if } 1 < i \leq |\bar{x}_{1}| \\
\bar{x}_{1}, \det_{i-|\bar{x}_{1}|}(\bar{x}_{2}, \dots, \bar{x}_{R}) & \text{if } i > |\bar{x}_{1}|
\end{aligned}$$

We define a replacement of the ith node with $y \in \mathcal{X}$ as the following function $\text{rep}_{i,y}$.

$$\operatorname{rep}_{i,y}(\bar{x}_{1},\ldots,\bar{x}_{R}) := \begin{cases} \bar{x}_{1},\ldots,\bar{x}_{R} & \text{if } i < 1\\ \operatorname{rep}_{i}(\bar{x}_{1},\ldots,\bar{x}_{R}) & \text{if } i = 1\\ \nu(\bar{x}_{1}) \left(\operatorname{rep}_{i-1,y}(\bar{\varrho}(\bar{x}_{1}))\right),\bar{x}_{2},\ldots,\bar{x}_{R} & \text{if } 1 < i \leq |\bar{x}_{1}|\\ \bar{x}_{1},\operatorname{rep}_{i-|\bar{x}_{1}|,y}(\bar{x}_{2},\ldots,\bar{x}_{R}) & \text{if } i > |\bar{x}_{1}| \end{cases}$$

Finally, we define an *insertion* of node $y \in \mathcal{X}$ as parent of the children l to r-1 of the ith node as the following function $ins_{i,y,l,r}$.

$$\operatorname{ins}_{i,y,l,r}(\bar{x}_1,\ldots,\bar{x}_R) := \begin{cases} \bar{x}_1,\ldots,\bar{x}_R & \text{if } i < 0 \\ \operatorname{ins}_{y,l,r}(\bar{x}_1,\ldots,\bar{x}_R) & \text{if } i = 0 \\ \nu(\bar{x}_1) \left(\operatorname{ins}_{i-1,y,l,r}(\bar{\varrho}(\bar{x}_1))\right) \bar{x}_2,\ldots,\bar{x}_R & \text{if } 1 \leq i \leq |\bar{x}_1| \\ \bar{x}_1,\operatorname{ins}_{i-|\bar{x}_1|,y,l,r}(\bar{x}_2,\ldots,\bar{x}_R) & \text{if } i > |\bar{x}_1| \end{cases}$$

We define the standard TED edit set $\Delta_{\mathcal{X}}$ for the alphabet \mathcal{X} as the following set: $\Delta_{\mathcal{X}} := \{ \operatorname{del}_i | i \in \mathbb{N} \} \cup \{ \operatorname{rep}_{i,y} | i \in \mathbb{N}, y \in \mathcal{X} \} \cup \{ \operatorname{ins}_{i,y,l,r} | i \in \mathbb{N}_0, l, r \in \mathbb{N}, y \in \mathcal{X} \}.$

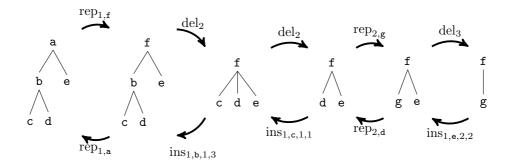


Figure 2: An illustration of a shortest edit script transforming the tree \bar{x} on the left to the tree \bar{y} on the right and a shortest edit script transforming the tree \bar{y} on the right to the tree \bar{x} on the left. The intermediate trees resulting from the application of single edits are shown in the middle.

Note that in all three cases we leave X unchanged if i < 1 or if i > |X| (except for insertions, where we only leave the forest unchanged if i < 0 or i > |X|).

In the remainder of this work, we focus on the TED according to the standard TED edit set $d_{\Delta_{\mathcal{X}}}$. An example of a shortest edit script according to these edits is shown in Figure 2. In particular, the script $\bar{\delta} = \operatorname{rep}_{1,\mathbf{f}}, \operatorname{del}_2, \operatorname{del}_2, \operatorname{rep}_{2,\mathbf{g}}, \operatorname{del}_3$ transforms the tree $\bar{x} = \mathbf{a}(\mathbf{b}(\mathbf{c},\mathbf{d}),\mathbf{e})$ into the tree $\bar{y} = \mathbf{f}(\mathbf{g})$. Because there is no shorter script to transform \bar{x} to \bar{y} , the tree edit distance between \bar{x} and \bar{y} is 5. Note that the deletion of \mathbf{b} in the tree $\mathbf{f}(\mathbf{b}(\mathbf{c},\mathbf{d}),\mathbf{e})$ results in the tree $\mathbf{f}(\mathbf{c},\mathbf{d},\mathbf{e})$, meaning that the children of \mathbf{b} , namely \mathbf{d} and \mathbf{e} , are now children of the parent of \mathbf{b} , namely \mathbf{f} . Deletions can also lead to trees becoming forests. In particular, a deletion of \mathbf{f} in the tree $\mathbf{f}(\mathbf{c},\mathbf{d},\mathbf{e})$ would result in the forest $\mathbf{c},\mathbf{d},\mathbf{e}$.

Conversely, an insertion takes (some of) the children of a tree and uses them as children of the newly inserted node. For example, the insertion of b in the tree f(c,d,e) in Figure 2 uses the children c and d of f as children for the new node b. Insertions can also be used to transform forests to trees by inserting a new node at the root. For example, the insertion $ins_{0,f,1,4}$ applied to the forest c,d,e would result in the tree f(c,d,e).

2.3 Cost Functions

Up until now we have defined the edit distance based on the length of the script required to transform \bar{x} into \bar{y} . However, we may want to regard some edits as more costly than others, because some elements may be easier to replace. This is reflected in manually defined cost matrices, such as the PAM and BLOSUM matrices from bioinformatics (S. Henikoff and J. G. Henikoff 1992). In general, we can express the cost of edits in terms of a cost function.

Definition 6 (Cost function). A cost function c over some alphabet \mathcal{X} with $-\notin \mathcal{X}$ is defined as a function $c:(\mathcal{X}\cup\{-\})\times(\mathcal{X}\cup\{-\})\to\mathbb{R}$, where - is called the special gap symbol.

Further, we define the cost of any edit δ in $\delta_{\mathcal{X}}$ as zero if $\delta(X) = X$, i.e. if the edit does not change the input forest, and otherwise as follows. We define the cost of a replacement $\operatorname{rep}_{i,y}$ with respect to some input forest X as $c(\operatorname{rep}_{i,y}, X) := c(x_i, y)$; we define the cost of a deletion del_i with respect to some input forest X as $c(\operatorname{del}_i, X) := c(x_i, -)$; and we define the cost of an insertion $\operatorname{ins}_{i,y,l,r}$ with respect to some input forest X as $c(\operatorname{ins}_{i,y,l,r}, X) := c(-,y)$.

Finally, we define the cost of an edit script $\bar{\delta} = \delta_1, \ldots, \delta_T$ with respect to some input forest X recursively as $c(\bar{\delta}, X) := c(\delta_1, X) + c((\delta_2, \ldots, \delta_T), \delta_1(X))$, with the base case $c(\epsilon, X) = 0$ for the empty script.

Intuitively, the cost of an edit script is just the sum over the costs of any single edit in the script.

As an example, consider our example script in Figure 2. For this script we obtain the cost:

$$\begin{split} c\Big(\mathrm{rep}_{1,\mathbf{f}}, \mathrm{del}_2, \mathrm{del}_2, \mathrm{rep}_{2,\mathbf{g}}, \mathrm{del}_3, \mathbf{a}(\mathbf{b}(\mathbf{c}, \mathbf{d}), \mathbf{e})\Big) \\ = &c(\mathbf{a}, \mathbf{f}) + c\Big(\mathrm{del}_2, \mathrm{del}_2, \mathrm{rep}_{2,\mathbf{g}}, \mathrm{del}_3, \mathbf{f}(\mathbf{b}(\mathbf{c}, \mathbf{d}), \mathbf{e})\Big) \\ = &c(\mathbf{a}, \mathbf{f}) + c(\mathbf{b}, -) + c\Big(\mathrm{del}_2, \mathrm{rep}_{2,\mathbf{g}}, \mathrm{del}_3, \mathbf{f}(\mathbf{c}, \mathbf{d}, \mathbf{e})\Big) \\ = &c(\mathbf{a}, \mathbf{f}) + c(\mathbf{b}, -) + c(\mathbf{c}, -) + c\Big(\mathrm{rep}_{2,\mathbf{g}}, \mathrm{del}_3, \mathbf{f}(\mathbf{d}, \mathbf{e})\Big) \\ = &c(\mathbf{a}, \mathbf{f}) + c(\mathbf{b}, -) + c(\mathbf{c}, -) + c(\mathbf{d}, \mathbf{g}) + c\Big(\mathrm{del}_3, \mathbf{f}(\mathbf{g}, \mathbf{e})\Big) \\ = &c(\mathbf{a}, \mathbf{f}) + c(\mathbf{b}, -) + c(\mathbf{c}, -) + c(\mathbf{d}, \mathbf{g}) + c(\mathbf{e}, -) + c\Big(\mathbf{c}, \mathbf{f}(\mathbf{g})\Big) \\ = &c(\mathbf{a}, \mathbf{f}) + c(\mathbf{b}, -) + c(\mathbf{c}, -) + c(\mathbf{d}, \mathbf{g}) + c(\mathbf{e}, -) + 0 \end{split}$$

Based on the notion of cost, we can generalize the TED as follows.

Definition 7 (Generalized Tree Edit Distance). Let \mathcal{X} be an alphabet, let $\Delta_{\mathcal{X}}$ be the standard TED edit set over \mathcal{X} , and let c be a cost function over \mathcal{X} . Then, the generalized TED over \mathcal{X} is defined as the function

$$d_c: \mathcal{T}(\mathcal{X}) \times \mathcal{T}(\mathcal{X}) \to \mathbb{R} \tag{5}$$

$$d_c(\bar{x}, \bar{y}) = \min_{\bar{\delta} \in \Delta_{\mathcal{X}}^*} \left\{ c(\bar{\delta}, \bar{x}) \middle| \bar{\delta}(\bar{x}) = \bar{y} \right\}$$
 (6)

As an example, consider the cost function c(x,y)=1 if $x\neq y$ and 0 if x=y. In that case, every edit (except for self-replacements) costs 1, such that the generalized edit distance again corresponds to the cost of the shortest edit script. If we change this cost function to be 0 for a replacement of a with f, our edit distance between the two example trees in Figure 2 decreases from 5 to 4. If we set the cost $c(\mathtt{a},\mathtt{a})=-1$, the edit distance becomes ill-defined, because we can always make an edit script cheaper by appending another self-replacement of a with a.

This begs the question: Which properties does the cost function c need to fulfill in order to ensure a "reasonable" edit distance? To answer this question, we first define what it means for a distance to be "reasonable". Here, we turn to the mathematical notion of a metric.

Definition 8 (Metric). Let \mathcal{X} be some set. A function $d: \mathcal{X} \times \mathcal{X} \to \mathbb{R}$ is called a *metric* if for all $x, y, z \in \mathcal{X}$ it holds:

$$\begin{array}{ll} d(x,y) \geq 0 & \text{(non-negativity)} \\ d(x,x) = 0 & \text{(self-equality)} \\ d(x,y) > 0 \text{ if } x \neq y & \text{(discernibility)} \\ d(x,y) = d(y,x) & \text{(symmetry)} \\ d(x,z) + d(z,y) \geq d(x,y) & \text{(triangular inequality)} \end{array}$$

All five of these properties make intuitive sense: We require a reasonable distance to not return negative values, we require that every object should have a distance of 0 to itself, we require that no two different objects can occupy the same space, we require that any object x is as far from y as y is from x, and we require that the fastest route from x to y is a straight line, that is, there is no point z through which we could travel such that we reach y faster from x compared to taking the direct distance.

Interestingly, it is relatively easy to show that the generalized TED is a metric if the cost function is a metric.

Theorem 1. If c is a metric over \mathcal{X} , then the generalized TED d_c is a metric over $\mathcal{T}(\mathcal{X})$. More specifically:

- 1. If c is non-negative, then d_c is non-negative.
- 2. If c is non-negative and self-equal, then d_c is self-equal.
- 3. If c is non-negative and discernible, then d_c is discernible.
- 4. If c is non-negative and symmetric, then d_c is symmetric.
- 5. If c is non-negative, d_c conforms to the triangular inequality.

Proof. Note that we require non-negativity as a pre-requisite for any of the metric conditions, because negative cost function values may lead to an ill-defined distance, as in the example above.

We now prove any of the four statements in turn:

Non-negativity: The TED is a sum of outputs of c. Because c is non-negative, d_c is as well.

Self-Equality: The empty edit script ϵ transforms \bar{x} to \bar{x} and has a cost of 0. Because d_c is non-negative, this is the cheapest edit sequence, therefore $d_c(\bar{x}, \bar{x}) = 0$ for all \bar{x} .

Discernibility: Let $\bar{x} \neq \bar{y}$ be two different trees and let $\bar{\delta} = \delta_1, \ldots, \delta_T$ be an edit script such that $\bar{\delta}(\bar{x}) = \bar{y}$. We now define $\bar{x}_0 = \bar{x}$ and \bar{x}_t recursively as $\delta_t(\bar{x}_{t-1})$ for all $t \in \{1, \ldots, T\}$. Accordingly, there must exist an $t \in \{1, \ldots, T\}$ such that $\bar{x}_t \neq \bar{x}_{t-1}$, otherwise $\bar{x} = \bar{y}$. However, in that case, the costs of δ_t must be c(x,y) for some $x \neq y$. Because c is discernible, c(x,y) > 0. Further, because c is non-negative, $c(\bar{\delta}, \bar{x})$ is a sum of non-negative contributions with at least one strictly positive contribution, which means that $c(\bar{\delta}, \bar{x}) > 0$. Because this reasoning applies for any script $\bar{\delta}$ with $\bar{\delta}(\bar{x}) = \bar{y}$, it holds: $d_c(\bar{x}, \bar{y}) > 0$.

Symmetry: Let $\bar{\delta} = \delta_1, \ldots, \delta_T$ be the cheapest edit script which transforms \bar{x} to \bar{y} . Now, we can inductively construct an *inverse* edit script as follows: If $\bar{\delta}$ is the empty script, then the empty script also transforms \bar{y} to \bar{x} . If $\bar{\delta}$ is not empty, consider the first edit δ_1 :

- If $\delta_1 = \operatorname{rep}_{i,y}$, we construct the edit $\delta_1^{-1} = \operatorname{rep}_{i,x_i}$. For this edit it holds: $\delta_1 \circ \delta_1^{-1}(\bar{x}) = \bar{x}$. Further, for the cost it holds: $c(\delta_1, \bar{x}) = c(x_i, y) = c(y, x_i) = c(\delta_1^{-1}, \delta(\bar{x}))$.
- If $\delta_1 = \operatorname{ins}_{i,y,l,r}$, we construct the edit $\delta_1^{-1} = \operatorname{del}_{i'}$ where i' is the index of the newly inserted node in the forest $\delta_1(\bar{x})$. Therefore, we obtain $\delta_1 \circ \delta_1^{-1}(\bar{x}) = \bar{x}$. Further, for the cost it holds: $c(\delta_1, \bar{x}) = c(-, y) = c(y, -) = c(\delta_1^{-1}, \delta_1(\bar{x}))$.
- If $\delta_1 = \text{del}_i$, we construct the edit $\delta_1^{-1} = \inf_{p_{\bar{x}}(i), x_i, r_{\bar{x}}(i), r_{\bar{x}}(i) + |\bar{\varrho}(\bar{x}_i)|}$. That is, we construct an insertion which re-inserts the node that has been deleted by δ_1 , and uses all its prior children. Therefore, we obtain $\delta_1 \circ \delta_1^{-1}(\bar{x}) = \bar{x}$. Further, for the cost it holds: $c(\delta_1, \bar{x}) = c(x_i, -) = c(-1, x_i) = c(\delta_1^{-1}, \delta(\bar{x}))$.

It follows by induction that we can construct an entire script $\bar{\delta}^{-1}$, which transforms \bar{y} to \bar{x} , because $\bar{x} = \bar{\delta} \circ \bar{\delta}^{-1}(\bar{x}) = \bar{\delta}^{-1}(\bar{y})$. Further, this script costs the same as $\bar{\delta}$, because $c(\bar{\delta}^{-1}, \bar{\delta}(\bar{x})) = c(\bar{\delta}^{-1}, \bar{y}) = c(\bar{\delta}, \bar{x})$.

Because $\bar{\delta}$ was by definition a cheapest edit script which transforms \bar{x} to \bar{y} we obtain: $d_c(\bar{y}, \bar{x}) \leq c(\bar{\delta}^{-1}, \bar{y}) = c(\bar{\delta}, \bar{x}) = d_c(\bar{x}, \bar{y})$. It remains to show that $d_c(\bar{y}, \bar{x}) \geq d_c(\bar{x}, \bar{y})$.

Assume that $d_c(\bar{y}, \bar{x}) < d_c(\bar{x}, \bar{y})$. Then, there is an edit script $\bar{\delta} = \delta_1, \dots, \delta_{T'}$ which transforms \bar{y} to \bar{x} and is cheaper than $d_c(\bar{x}, \bar{y})$. However, using the same argument as before, we can generate an inverse edit script $\bar{\delta}^{-1}$ with the same cost as $\bar{\delta}$ that transforms \bar{x} to \bar{y} , such that $d_c(\bar{x}, \bar{y}) \leq d_c(\bar{y}, \bar{x}) < d_c(\bar{x}, \bar{y})$, which is a contradiction. Therefore $d_c(\bar{y}, \bar{x}) = d_c(\bar{x}, \bar{y})$.

Triangular Inequality: Assume that there are three trees \bar{x} , \bar{y} , and \bar{z} , such that $d_c(\bar{x}, \bar{z}) + d_c(\bar{z}, \bar{y}) < d_c(\bar{x}, \bar{y})$. Now, let $\bar{\delta}$ and $\bar{\delta}'$ be cheapest edit scripts which transform \bar{x} to \bar{z} and \bar{z} to \bar{y} respectively, that is, $\bar{\delta}(\bar{x}) = \bar{z}$, $\bar{\delta}'(\bar{z}) = \bar{y}$, $c(\bar{\delta}, \bar{x}) = d_c(\bar{x}, \bar{z})$, and $c(\bar{\delta}', \bar{z}) = d_c(\bar{z}, \bar{y})$. The concatenation of both scripts $\bar{\delta}'' = \bar{\delta} \oplus \bar{\delta}'$ is per construction a script such that $\bar{\delta}''(\bar{x}) = \bar{\delta} \oplus \bar{\delta}'(\bar{x}) = \bar{y}$ and $c(\bar{\delta}'', \bar{x}) = c(\bar{\delta}, \bar{x}) + c(\bar{\delta}', \bar{z}) = d_c(\bar{x}, \bar{z}) + d_c(\bar{z}, \bar{y})$. It follows that $d_c(\bar{x}, \bar{z}) + d_c(\bar{z}, \bar{y}) < d_c(\bar{x}, \bar{y}) \leq d_c(\bar{x}, \bar{z}) + d_c(\bar{z}, \bar{y})$ which is a contradiction. Therefore, the triangular inequality holds.

As an example of the symmetry part of the proof, consider again Figure 2. Here, the inverse script for $\bar{\delta} = \text{rep}_{1,f}, \text{del}_2, \text{del}_2, \text{rep}_{2,g}, \text{del}_3$ is $\bar{\delta}^{-1} = \text{ins}_{1,e,2,2}, \text{rep}_{2,d}, \text{ins}_{1,c,1,1}, \text{ins}_{1,b,1,3}, \text{rep}_{1,a}$. For the cost we obtain:

$$\begin{split} c\Big(&\inf_{1,\mathsf{e},2,2}, \operatorname{rep}_{2,\mathsf{d}}, \inf_{1,\mathsf{c},1,1}, \inf_{1,\mathsf{b},1,3}, \operatorname{rep}_{1,\mathsf{a}}, \mathsf{f}(\mathsf{g})\Big) \\ =&c(-,\mathsf{e}) + c\Big(\operatorname{rep}_{2,\mathsf{d}}, \inf_{1,\mathsf{c},1,1}, \inf_{1,\mathsf{b},1,3}, \operatorname{rep}_{1,\mathsf{a}}, \mathsf{f}(\mathsf{g},\mathsf{e})\Big) \\ =&c(-,\mathsf{e}) + c(\mathsf{g},\mathsf{d}) + c\Big(\inf_{1,\mathsf{c},1,1}, \inf_{1,\mathsf{b},1,3}, \operatorname{rep}_{1,\mathsf{a}}, \mathsf{f}(\mathsf{d},\mathsf{e})\Big) \\ =&c(-,\mathsf{e}) + c(-,\mathsf{d}) + c(-,\mathsf{c}) + c\Big(\inf_{1,\mathsf{b},1,3}, \operatorname{rep}_{1,\mathsf{a}}, \mathsf{f}(\mathsf{c},\mathsf{d},\mathsf{e})\Big) \\ =&c(-,\mathsf{e}) + c(-,\mathsf{d}) + c(-,\mathsf{c}) + c(-,\mathsf{b}) + c\Big(\operatorname{rep}_{1,\mathsf{a}}, \mathsf{f}(\mathsf{b}(\mathsf{c},\mathsf{d}),\mathsf{e})\Big) \\ =&c(-,\mathsf{e}) + c(-,\mathsf{d}) + c(-,\mathsf{c}) + c(\mathsf{g},\mathsf{b}) + c(\mathsf{f},\mathsf{a}) + c\Big(\epsilon, \mathsf{a}(\mathsf{b}(\mathsf{c},\mathsf{d}),\mathsf{e})\Big) \\ =&c(-,\mathsf{e}) + c(-,\mathsf{d}) + c(-,\mathsf{c}) + c(\mathsf{g},\mathsf{b}) + c(\mathsf{f},\mathsf{a}) + 0 \end{split}$$

2.4 Mappings

While edit scripts capture the intuitive notion of editing a tree, they are not a viable representation to develop an efficient algorithm. In particular, edit scripts are highly redundant, in the sense that there may be many different edit scripts which transform a tree \bar{x} into a tree \bar{y} and have the same cost. For example, to transform the tree $\bar{x} = \mathbf{a}(\mathbf{b}(\mathbf{c}, \mathbf{d}), \mathbf{e})$ to the tree $\bar{y} = \mathbf{f}(\mathbf{g})$, we can not only use the edit script in Figure 2, but we could also use the script del₅, del₃, del₂, rep_{2,g}, rep_{1,f}, which has the same cost, irrespective of the cost function. To avoid these redundancies, we need a representation which is invariant against changes in order of the edits, and instead just counts which nodes are replaced, which nodes are deleted, and which nodes are inserted. This representation is offered by the concept of a tree mapping.

Definition 9 (Tree Mapping). Let X and Y be forests over some alphabet \mathcal{X} , and let m = |X| and n = |Y|.

A tree mapping between X and Y is defined as a set of tuples $M \subseteq \{1, ..., m\} \times \{1, ..., n\}$, such that for all $(i, j), (i', j') \in M$ it holds:

- 1. Each node of X is assigned to at most one node in Y, i.e. $i = i' \Rightarrow j = j'$.
- 2. Each node of Y is assigned to at most one node in X, i.e. $j = j' \Rightarrow i = i'$.
- 3. The mapping preserves the pre-order of both trees, i.e. $i \ge i' \iff j \ge j'$.
- 4. The mapping preserves the ancestral ordering in both trees, that is: if the subtree rooted at i is an ancestor of the subtree rooted at i', then the subtree rooted at j is also an ancestor of the subtree rooted at j', and vice versa.

The four constraints in the definition of a tree mapping have the purpose to ensure that we can find a corresponding edit script for each mapping. As an example, consider again our two trees $\bar{x} = \mathbf{a}(\mathbf{b}(\mathbf{c}, \mathbf{d}), \mathbf{e})$ and $\bar{y} = \mathbf{f}(\mathbf{g})$ from Figure 2. The mapping corresponding to the edit script in this figure would be $M = \{(1,1), (4,2)\}$ because node $x_1 = \mathbf{a}$ is replaced with node $y_1 = \mathbf{f}$ and node $x_4 = \mathbf{d}$ is replaced with node $y_2 = \mathbf{g}$. All remaining nodes are deleted and inserted respectively. The mapping $M = \{\}$ would correspond to deleting all nodes in \bar{x} and then inserting all nodes in \bar{y} , which is a valid mapping but would be more costly.

The set $M = \{(1,1), (1,2)\}$ would not be a valid mapping because the node $x_1 = \mathbf{a}$ is assigned to multiple nodes in \bar{y} and thus we can not construct an edit script corresponding to this mapping. For such an edit script we would need a "copy" edit. For the same reason, the set $M = \{(1,1), (2,1)\}$ is not a valid mapping. Here, the node $y_1 = \mathbf{f}$ is assigned to multiple nodes in \bar{x} .

 $M = \{(1,2),(2,1)\}$ is an example of a set that is not a valid tree mapping because of the third criterion. To construct an edit script corresponding to this mapping we would need a "swap" edit, i.e. an edit which can exchange nodes $x_1 = \mathbf{a}$ and $x_2 = \mathbf{b}$ in \bar{x} . Finally, the set $M = \{(3,1),(5,2)\}$ is not a valid mapping due to the fourth criterion. In particular, the subtree $\bar{y}_1 = \mathbf{f}(\mathbf{g})$ is an ancestor of

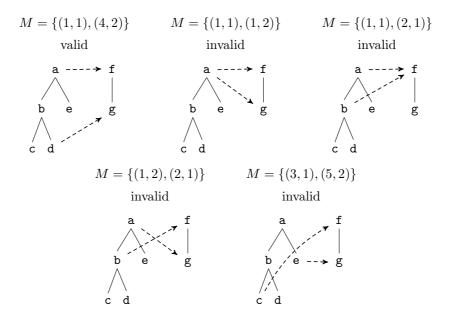


Figure 3: One example mapping between the trees $\bar{x} = a(b(c,d),e)$ and $\bar{y} = f(g)$ (top left) and four sets which are not valid mappings due to violations of one of the four mapping constraints.

the subtree $\bar{y}_2 = g$ in \bar{y} , but the subtree $\bar{x}_3 = c$ is *not* an ancestor of the subtree $\bar{x}_5 = e$. This last criterion is more subtle, but you will find that each edit we can apply - be it replacement, deletion, or insertion - preserves the ancestral order in the tree. Conversely, this means that we can not make a node an ancestor of another node if it was not before. This also makes intuitive sense because it means that nodes can not be mapped to nodes in completely distinct subtrees.

Now that we have considered some examples, it remains to show that we can construct a corresponding edit script for each mapping in general.

Theorem 2. Let X and Y be forests over some alphabet \mathcal{X} and let M be a tree mapping between X and Y. Then, the output of Algorithm 1 for X, Y and M is an edit script $\bar{\delta}_M$ using only replacements, deletions, and insertions, such that $\bar{\delta}_M(X) = Y$ and

- 1. If $(i,j) \in M$, then the edit rep_{i,y_i} is part of the script.
- 2. For all i which are not part of the mapping i.e. $\nexists j:(i,j)\in M$ the edit del_i is part of the script.
- 3. For all j which are not part of the mapping i.e. $\nexists i:(i,j)\in M$ an edit $\operatorname{ins}_{p(j),y_j,r_j,r_j+R_j}$ for some R_j is part of the script.

Further, no other edits are part of $\bar{\delta}_M$ than the edits mentioned above. Algorithm 1 runs in $\mathcal{O}(m+n)$ worst-case time.

Proof. The three constraints are fulfilled because we iterate over all entries (i,j) and create one replacement per such entry, we iterate over all $i \in I$ and create one deletion per such i, and we iterate over all $j \in J$ and we create one insertion per such entry. It is also clear that $\mathcal{O}(m+n)$ because we iterate over all $i \in \{1, ..., m\}$ and over all $j \in \{1, ..., n\}$. Assuming that I and J permit insertion as well as containment tests in constant time, and that the list concatenations in num-descendants are possible in constant time, this leaves us with $\mathcal{O}(m+n)$.

It is less obvious that $\delta_M(X) = Y$. We show this by an induction over the cardinality of M. First, consider $M = \emptyset$. In that case, we obtain $I = \{1, \ldots, m\}$, $J = \{1, \ldots, n\}$, and $R_0 = \ldots = R_n = 0$. Therefore, the resulting script is $\bar{\delta}_M = \text{del}_m, \ldots, \text{del}_1, \text{ins}_{p(1), \nu(\bar{y}_1), r_1, r_1}, \ldots, \text{ins}_{p(n), \nu(\bar{y}_n), r_n, r_n}$. This script obviously first deletes all nodes in X and then inserts all nodes from Y in the correct configuration.

Now assume that the theorem holds for all mappings M between X and Y with $|M| \leq k$, and consider a mapping M between X and Y with |M| = k + 1.

Let (i, j) be the entry of M with smallest j, let $M' = M \setminus \{(i, j)\}$ and let $\delta_{M'}$ be the corresponding edit script for M' according to Algorithm 1. Per induction, $\delta_{M'}(X) = Y$.

Now, let $I = \{i | \nexists j : (i,j) \in M\}$, $I' = \{i | \nexists j : (i,j) \in M'\}$, $J = \{j | \nexists i : (i,j) \in M\}$, and $J' = \{j | \nexists i : (i,j) \in M'\}$. We observe that $I' = I \cup \{i\}$ and $J' = J \cup \{j\}$, so our resulting script $\bar{\delta}_M$ will not delete node i and not insert node j, but otherwise contain all deletions and insertions of script $\bar{\delta}_{M'}$. We also know that node x_i will be replaced with node y_j , such that all nodes of Y are contained after applying $\bar{\delta}_M$. It remains to show that node y_j is positioned correctly in $\bar{\delta}_M(X)$, such that $\bar{\delta}_M(X) = Y$.

Let $P_J(j)$ be a set that is recursively defined as $P_J(j') = \emptyset$ if $j' \notin J$, and $P_J(j') = \{j'\} \cup P_J(p_{\bar{y}}(j'))$ if $j' \in J$. In other words, $P_J(j')$ contains all ancestors of j', until we find an ancestor that is not inserted. Now, consider all inserted ancestors of j, that is, $P_J(p_{\bar{y}}(j))$. Further, let $(n, R_0, \ldots, R_n) =$ num-descendants(Y, 0, J) and $(n, R'_0, \ldots, R'_n) =$ num-descendants(Y, 0, J'). For all elements $j' \in P_J(p_{\bar{y}}(j))$ we obtain $R_{j'} = R'_{j'} + 1$, and for all other nodes we obtain $R_{j'} = R'_{j'}$. In other words, all ancestors of j which are inserted use one more child compared to before, but no other node will. This additional child is j, such that the ancestral structure is preserved and we obtain $\bar{\delta}_M(X) = Y$.

Algorithm 1 An algorithm to transform a mapping M into a corresponding edit script $\bar{\delta}_M$ according to Theorem 2.

```
function MAP-TO-SCRIPT (Two forests X and Y, a tree mapping M between X and Y.)
     I \leftarrow \{i | \nexists j : (i,j) \in M\}.
     J \leftarrow \{j | \nexists i : (i,j) \in M\}.
     Initialize \bar{\delta} as empty.
     for (i,j) \in M do
           \bar{\delta} \leftarrow \bar{\delta} \oplus \operatorname{rep}_{i,y_j}.
                                                                                                                                      ▷ replacements
     end for
     for i \in I in descending order do
           \bar{\delta} \leftarrow \bar{\delta} \oplus \operatorname{del}_i.
                                                                                                                                            ▶ deletions
     end for
     (n, R_0, \ldots, R_n) \leftarrow \text{NUM-DESCENDANTS}(Y, 0, J).
                                                                                         > the number of children for each inserted
node
     for j \in J in ascending order do
           \bar{\delta} \leftarrow \bar{\delta} \oplus \operatorname{ins}_{p(j),\nu(\bar{y}_j),r_j,r_j+R_j},
                                                                                                                                           \triangleright insertions
     end for
     return \bar{\delta}.
end function
function NUM-DESCENDANTS (A forest Y = \bar{y}_1, \dots, \bar{y}_R, an index j, and an index set J)
     R \leftarrow \epsilon.
     \tilde{R} \leftarrow 0.
                                                                           ▶ The number of mapped descendants of this forest
     for r \leftarrow 1, \dots, R do
          j \leftarrow j+1.
           (j', \tilde{R}_j, \dots, \tilde{R}_{j'}) \leftarrow \text{NUM-DESCENDANTS}(\bar{\varrho}(\bar{y}_r), j, J).
           \bar{R} \leftarrow \bar{R} \oplus \tilde{R}_i, \dots, \tilde{R}_{i'}.
          if j \notin J then
                R \leftarrow R + 1.
           else
                \tilde{R} \leftarrow \tilde{R} + \tilde{R}_i.
           end if
           j \leftarrow j'.
     end for
     return (j, \tilde{R} \oplus \bar{R}).
end function
```

As an example, consider again the mapping $M = \{(1,1),(4,2)\}$ between the trees $\bar{x} = \mathbf{a}(\mathbf{b}(\mathbf{c},\mathbf{d}),\mathbf{e})$

and $\bar{y} = f(g)$ from Figure 2. Here we have the non-mapped nodes $I = \{2, 3, 5\}$ and $J = \{\}$. Therefore, Algorithm 1 returns the script $\text{rep}_{1,f}$, $\text{rep}_{4,g}$, del_5 , del_3 , del_2 . Note that deletions are done in descending order to ensure that the pre-order indices in the tree do not change for intermediate trees.

For the inverse mapping $M = \{(1,1),(2,4)\}$ between \bar{y} and \bar{x} we have $I = \{\}$ and $J = \{2,3,5\}$. Further, the output of num-descendants is $(n = 5, R_0 = 1, R_1 = 1, R_2 = 1, R_3 = 0, R_4 = 0, R_5 = 0)$. Therefore, we obtain the script $\text{rep}_{1,a}$, $\text{rep}_{2,d}$, $\text{ins}_{1,b,1,2}$, $\text{ins}_{2,c,1,1}$, $\text{ins}_{1,e,2,2}$.

Our next task is to demonstrate that the inverse direction is also possible, that is, we can find a corresponding mapping for each script.

Theorem 3. Let X and Y be forests over some alphabet \mathcal{X} , and let $\bar{\delta}$ be an edit script such that $\bar{\delta}(X) = Y$. Then, the following, recursively defined set $M_{\bar{\delta}}$, is a mapping between X and Y:

$$M_{\epsilon} := \{(1,1), \dots, (m,m)\}$$

$$M_{\delta_{1},\dots,\delta_{T}-1} \qquad if \ \delta_{T} = \operatorname{rep}_{j,y_{j}}$$

$$\{(i,j')|(i,j') \in M_{\delta_{1},\dots,\delta_{T-1}}, j' < j\} \quad \cup \quad if \ \delta_{T} = \operatorname{del}_{j}$$

$$\{(i,j'-1)|(i,j') \in M_{\delta_{1},\dots,\delta_{T-1}}, j' > j\} \quad \cup \quad if \ \delta_{T} = \operatorname{del}_{j}$$

$$\{(i,j')|(i,j') \in M_{\delta_{1},\dots,\delta_{T-1}}, j' < j\} \quad \cup \quad if \ \delta_{T} = \operatorname{ins}_{p(j),y_{j},r_{j},r_{j}+R_{j}}$$

where R_j is the number of children of \bar{x}_j .

Proof. We prove the claim via induction over the length of $\bar{\delta}$. M_{ϵ} obviously conforms to all mapping constraints.

Now, assume that the claim is true for all scripts $\bar{\delta}$ with $|\bar{\delta}| \leq T$ and consider a script $\bar{\delta} = \delta_1, \ldots, \delta_{T+1}$. Let $\bar{\delta}' = \delta_1, \ldots, \delta_T$. Due to induction, we know that $M_{\bar{\delta}'}$ is a valid mapping between X and $\bar{\delta}'(X)$. Now, consider the last edit δ_{T+1} .

First, we observe that, if $M_{\bar{\delta}'}$ fulfills the first three criteria of a mapping, $M_{\bar{\delta}}$ does as well, because we never introduce many-to-one mappings and respect the pre-order. The only criterion left in question is the fourth, namely whether $M_{\bar{\delta}}$ respects the ancestral ordering of Y.

If δ_{T+1} is a replacement, the tree structure of $\delta(X)$ is the same as for $\bar{\delta}'(X)$. Therefore, $M_{\bar{\delta}} = M_{\bar{\delta}'}$ is also a valid mapping between X and Y.

If δ_{T+1} is a deletion del_j , then node y_j in $\bar{\delta}'(X)$ is missing from Y and all subtrees with pre-order indices higher than j decrease their index by one, which is reflected by $M_{\bar{\delta}}$. Further, $M_{\bar{\delta}}$ only removes a tuple, but does not add a tuple, such that all ancestral relationships present in $M_{\bar{\delta}}$ were also present in $M_{\bar{\delta}'}$. Finally, a deletion does not break any of the ancestral relationships because any ancestor of \bar{y}_j remains an ancestor of all children of \bar{y}_j in Y. Therefore, $M_{\bar{\delta}}$ is a valid mapping between X and Y.

If δ_{T+1} is an insertion $\inf_{p(j),y_j,r_j,r_j+R_j}$, then y_j is a new node in Y and all subtrees with preorder indices as high or higher than j in $\bar{\delta}'(X)$ increase their index by one, which is reflected by $M_{\bar{\delta}}$. Further, $M_{\bar{\delta}}$ leaves all tuples intact, such that all ancestral relationships of $M_{\bar{\delta}'}$ are preserved. Finally, an insertion does not break any ancestral relationships because $\bar{y}_{p(j)}$ is still an ancestor of all nodes it was before, except that there is now a new node y_j in between. Therefore, $M_{\bar{\delta}}$ is a valid mapping between X and Y.

As an example, consider the edit scripts shown in Figure 2. For the script $\bar{\delta} = \text{rep}_{1,f}$, del_2 , del_2 , $\text{rep}_{2,g}$, del_3 , which transforms the tree $\bar{x} = \mathsf{a}(\mathsf{b}(\mathsf{c},\mathsf{d}),\mathsf{e})$ into the tree $\bar{y} = \mathsf{f}(\mathsf{g})$, we obtain the following mappings M_t after the tth edit:

$$\begin{split} M_0 &= \{(1,1),(2,2),(3,3),(4,4),(5,5)\} & \text{initial} \\ M_1 &= \{(1,1),(2,2),(3,3),(4,4),(5,5)\} & \text{rep}_{1,\mathbf{f}}, \text{del}_2 \\ M_2 &= \{(1,1),(3,2),(4,3),(5,4)\} & \text{rep}_{1,\mathbf{f}}, \text{del}_2 \\ M_3 &= \{(1,1),(4,2),(5,3)\} & \text{rep}_{1,\mathbf{f}}, \text{del}_2, \text{del}_2 \\ M_4 &= \{(1,1),(4,2),(5,3)\} & \text{rep}_{1,\mathbf{f}}, \text{del}_2, \text{del}_2, \text{rep}_{2,\mathbf{g}} \\ M_5 &= \{(1,1),(4,2)\} & \text{rep}_{1,\mathbf{f}}, \text{del}_2, \text{del}_2, \text{rep}_{2,\mathbf{g}}, \text{del}_3 \end{split}$$

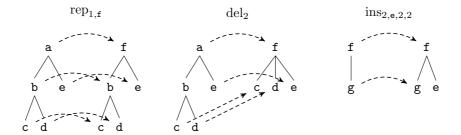


Figure 4: An illustration of the recursive construction of the corresponding mapping $M_{\bar{\delta}}$ for three edits from Figure 2. For each edit of the script, the mapping is updated to be consistent with all edits up until now. In particular, the mapping starts as a one-to-one mapping, is left unchanged for all replacements, and is shifted for all deletions and insertions

Conversely, for the script $\bar{\delta}^{-1} = ins_{1,e,2,2}, rep_{2,d}, ins_{1,c,1,1}, ins_{1,b,1,3}, rep_{1,a}$, which transforms \bar{y} into \bar{x} , we obtain the following mappings.

$$\begin{split} M_0 &= \{(1,1),(2,2)\} & \text{initial} \\ M_1 &= \{(1,1),(2,2)\} & \text{ins}_{1,\mathbf{e},2,2} \\ M_2 &= \{(1,1),(2,2)\} & \text{ins}_{1,\mathbf{e},2,2}, \text{rep}_{2,\mathbf{d}} \\ M_3 &= \{(1,1),(2,3)\} & \text{ins}_{1,\mathbf{e},2,2}, \text{rep}_{2,\mathbf{d}}, \text{ins}_{1,\mathbf{c},1,1} \\ M_4 &= \{(1,1),(2,4)\} & \text{ins}_{1,\mathbf{e},2,2}, \text{rep}_{2,\mathbf{d}}, \text{ins}_{1,\mathbf{c},1,1}, \text{ins}_{1,\mathbf{b},1,3} \\ M_5 &= \{(1,1),(2,4)\} & \text{ins}_{1,\mathbf{e},2,2}, \text{rep}_{2,\mathbf{d}}, \text{ins}_{1,\mathbf{c},1,1}, \text{ins}_{1,\mathbf{b},1,3}, \text{rep}_{1,\mathbf{a}} \\ \end{split}$$

The influence of the different kinds of edits on the mapping is also illustrated in Figure 4.

Now that we have shown that edit scripts and mappings can be related on a structural level, it remains to show that they are also related in terms of *cost*. To that end, we need to define the cost of a mapping:

Definition 10 (Mapping cost). Let X and Y be forests over some alphabet \mathcal{X} , and let c be a cost function over \mathcal{X} . Further, let M be a mapping between X and Y, let $I = \{i \in \{1, ..., |X|\} | \nexists j : (i, j) \in M\}$, and let $J = \{j \in \{1, ..., |Y|\} | \nexists i : (i, j) \in M\}$

The cost of the mapping M is defined as:

$$c(M, X, Y) = \sum_{(i,j)\in M} c(x_i, y_j) + \sum_{i\in I} c(x_i, -) + \sum_{j\in J} c(-, y_j)$$
(7)

For example, consider the mapping $M = \{(1,1),(4,2)\}$ between the trees in Figure 2. This mapping has cost

$$c\Big(\{(1,1),(4,2)\},\mathtt{a}(\mathtt{b}(\mathtt{c},\mathtt{d}),\mathtt{e}),\mathtt{f}(\mathtt{g})\Big) = c(\mathtt{a},\mathtt{f}) + c(\mathtt{d},\mathtt{g}) + c(\mathtt{b},-) + c(\mathtt{c},-) + c(\mathtt{e},-)$$

Note that this is equivalent to the cost of the edit script $\bar{\delta} = \operatorname{rep}_{1,\mathbf{f}}, \operatorname{del}_2, \operatorname{del}_2, \operatorname{rep}_{2,\mathbf{g}}, \operatorname{del}_3$. However, the cost of an edit script is not always equal to the cost of its corresponding mapping. In general, we may have corresponding mappings which are cheaper then the scripts. For example, consider the two trees $\bar{x} = \mathbf{a}$ and $\bar{y} = \mathbf{b}$ and the script $\bar{\delta} = \operatorname{rep}_{1,\mathbf{c}}, \operatorname{rep}_{1,\mathbf{b}}$, which transforms \bar{x} to \bar{y} . Here, the corresponding mapping is $M = \{(1,1)\}$ with the cost $c(M,\bar{x},\bar{y}) = c(\mathbf{a},\mathbf{b})$. However, the cost of the edit script is $c(\bar{\delta},\bar{x}) = c(\mathbf{a},\mathbf{c}) + c(\mathbf{c},\mathbf{b})$, which will be at least as expensive if the cost function conforms to the triangular inequality.

In general, we obtain the following result:

Theorem 4. Let X and Y be forests over some alphabet \mathcal{X} and c be a cost function over \mathcal{X} . Further, let $\bar{\delta}$ be an edit script with $\bar{\delta}(X) = Y$, and let M be a mapping between X and Y. Then it holds:

- 1. The corresponding script $\bar{\delta}_M$ for M according to Algorithm 1 has the same cost as M, that is: $c(M, X, Y) = c(\bar{\delta}_M, X)$.
- 2. If c is non-negative, self-equal, and conforms to the triangular inequality, the corresponding mapping $M_{\bar{\delta}}$ for $\bar{\delta}$ according to Theorem 3 is at most as expensive as $\bar{\delta}$, that is: $c(M_{\bar{\delta}}, X, Y) \leq c(\bar{\delta}, X)$.

Proof. Let m = |X| and n = |Y|, let $I = \{i \in \{1, ..., m\} | \nexists j : (i, j) \in M\}$, and let $J = \{j \in \{1, ..., n\} | \nexists i : (i, j) \in M\}$

1. Due to Theorem 2 we know that the script $\bar{\delta}_M$ for M contains exactly one replacement $\operatorname{rep}_{i,y_j}$ per entry $(i,j) \in M$, exactly one deletion del_i per unmapped index $i \in I$ and exactly one insertion $\operatorname{ins}_{p(j),y_j,r_j,r+R_j}$ per unmapped index $j \in J$. Therefore, the cost of $\bar{\delta}_M$ is:

$$c(\bar{\delta}_M, X) = \sum_{(i,j) \in M} c(x_i, y_j) + \sum_{i \in I} c(x_i, -) + \sum_{j \in J} c(-, y_j)$$
(8)

which is per definition equal to c(M, X, Y).

2. We show this claim via induction over the length of $\bar{\delta}$. First, consider the case $\bar{\delta} = \epsilon$. Then, X = Y and $M_{\bar{\delta}} = \{(1,1), \ldots, (m,m)\}$. Because c is self-equal, we obtain for the cost of $M_{\bar{\delta}}$:

$$c(M_{\bar{\delta}}, X, Y) = \sum_{(i,j) \in M} c(x_i, y_j) = \sum_{i=1}^{m} c(x_i, x_i) = 0 = c(\epsilon, X)$$

Now, assume that the claim holds for all $\bar{\delta}'$ with $|\bar{\delta}'| \leq T$, and consider a script $\bar{\delta} = \delta_1, \ldots, \delta_{T+1}$. Let $\bar{\delta}' = \delta_1, \ldots, \delta_T$ and let $Y' = \bar{\delta}'(X)$. Then, we consider the last edit δ_{T+1} and distinguish the following cases:

If $\delta_{T+1} = \operatorname{rep}_{j,y}$ we have $M_{\bar{\delta}} = M_{\bar{\delta}'}$. Further, if there is an $i \in \{1, \dots, m\}$ such that $(i, j) \in M_{\bar{\delta}'}$, we obtain for the cost:

$$c(\bar{\delta}, X) = c(\bar{\delta}', X) + c(y'_j, y_j) \overset{\text{Induction}}{\geq} c(M_{\bar{\delta}'}, X, Y') + c(y'_j, y_j)$$

$$= c(M_{\bar{\delta}}, X, Y) - c(x_i, y_j) + c(x_i, y'_j) + c(y'_j, y_j)$$

$$\geq c(M_{\bar{\delta}}, X, Y) - c(x_i, y_j) + c(x_i, y_j) = c(M_{\bar{\delta}}, X, Y)$$

In case there is no $i \in \{1, ..., m\}$ such that $(i, j) \in M_{\bar{\delta}'}$, we obtain for the cost:

$$c(\bar{\delta}, X) = c(\bar{\delta}', X) + c(y'_{j}, y_{j}) \stackrel{\text{Induction}}{\geq} c(M_{\bar{\delta}'}, X, Y') + c(y'_{j}, y_{j})$$

$$= c(M_{\bar{\delta}}, X, Y) - c(-, y_{j}) + c(-, y'_{j}) + c(y'_{j}, y_{j})$$

$$\geq c(M_{\bar{\delta}}, X, Y) - c(-, y_{j}) + c(-, y_{j}) = c(M_{\bar{\delta}}, X, Y)$$

If $\delta_{T+1} = \operatorname{del}_j$, consider first the case that there exists some i such that $(i, j) \in M_{\bar{\delta}'}$. Then, we obtain for the cost:

$$c(\bar{\delta}, X) = c(\bar{\delta}', X) + c(y_j', -) \stackrel{\text{Induction}}{\geq} c(M_{\bar{\delta}'}, X, Y') + c(y_j', -)$$

$$= c(M_{\bar{\delta}}, X, Y) - c(x_i, -) + c(x_i, y_j') + c(y_j', -)$$

$$\geq c(M_{\bar{\delta}}, X, Y) - c(x_i, -) + c(x_i, -) = c(M_{\bar{\delta}}, X, Y)$$

If there exists no such i, we obtain for the cost:

$$c(\bar{\delta}, X) = c(\bar{\delta}', X) + c(y_j', -) \overset{\text{Induction}}{\geq} c(M_{\bar{\delta}'}, X, Y') + c(y_j', -)$$
$$= c(M_{\bar{\delta}}, X, Y) + c(-, y_j') + c(y_j', -) \geq c(M_{\bar{\delta}}, X, Y)$$

Finally, if $\delta_{T+1} = \inf_{p(j), y_i, r_i, r_i + R_i}$, we obtain for the cost

$$\begin{split} c(\bar{\delta},X) &= c(\bar{\delta}',X) + c(-,y_j) \overset{\text{Induction}}{\geq} c(M_{\bar{\delta}'},X,Y') + c(-,y_j) \\ &= c(M_{\bar{\delta}},X,Y) - c(-,y_j) + c(-,y_j) = c(M_{\bar{\delta}},X,Y) \end{split}$$

This concludes our proof by induction.

It follows directly that we can compute the TED by computing the cheapest mapping instead of the cheapest edit script.

Theorem 5. Let X and Y be forests over some alphabet \mathcal{X} and c be a cost function over \mathcal{X} that is non-negative, self-equal, and conforms to the triangular inequality. Then it holds:

$$\min_{\bar{\delta} \in \Delta_{\mathcal{X}}^*} \{ c(\bar{\delta}, X) | \bar{\delta}(X) = Y \} = \\ \min\{ c(M, X, Y) | M \text{ is a tree mapping between } X \text{ and } Y \}$$
 (9)

Proof. First, we define two abbreviations for the minima, namely:

$$\begin{split} d_c^{\text{script}}(X,Y) &:= \min_{\bar{\delta} \in \Delta_{\mathcal{X}}^*} \{c(\bar{\delta},X) | \bar{\delta}(X) = Y\} \\ d_c^{\text{map}}(X,Y) &:= \min\{c(M,X,Y) | M \text{ is a tree mapping between } X \text{ and } Y\} \end{split}$$

Let $\bar{\delta}$ be an edit script such that $\bar{\delta}(X)=Y$ and $c(\bar{\delta},X)=d_c^{\rm script}(X,Y)$. Then, we know due to Theorem 4 that the corresponding mapping $M_{\bar{\delta}}$ is at most as expensive as $\bar{\delta}$, i.e. $c(M_{\bar{\delta}},X,Y) \leq c(\bar{\delta},X)=d_c^{\rm script}(X,Y)$. This implies: $d_c^{\rm map}(X,Y) \leq d_c^{\rm script}(X,Y)$.

Conversely, let M be a tree mapping between X and Y, such that $c(M,X,Y)=d_c^{\rm map}(X,Y)$.

Conversely, let M be a tree mapping between X and Y, such that $c(M, X, Y) = d_c^{\text{map}}(X, Y)$. Then, we know due to Theorem 4 that the corresponding edit script $\bar{\delta}_M$ has the same cost as M, i.e. $c(\bar{\delta}_M, X) = c(M, X, Y)$. This implies: $d_c^{\text{script}}(X, Y) \leq d_c^{\text{map}}(X, Y)$.

This concludes our theory on edit scripts, cost functions, and mappings. We have now laid enough groundwork to efficiently compute the TED.

3 The Dynamic Programming Algorithm

To compute the TED between two trees \bar{x} and \bar{y} efficiently, we require a way to decompose the TED into parts, such that we can compute the distance between subtrees of \bar{x} and \bar{y} and combine those partial TEDs to an overall TED. In order to do that, we need to define what we mean by "partial trees".

Definition 11 (subforest). Let X be a forest of size m = |X|. Further, let $i, j \in \{1, ..., m\}$ with $i \leq j$. We define the *subforest* of X from i to j, denoted as X[i, j], as the first output of Algorithm 2 for the input X, i, j, and 0.

As an example, consider the left tree $\bar{x} = \mathtt{a}(\mathtt{b}(\mathtt{c},\mathtt{d}),\mathtt{e})$ from Figure 1 (left). For this example, we find: $X[1,1] = \mathtt{a}, X[2,4] = \mathtt{b}(\mathtt{c},\mathtt{d}), X[3,5] = \mathtt{c},\mathtt{d},\mathtt{e},$ and $X[2,1] = \epsilon.$

Note that $X[2,4] = \bar{x}_2$, that is: The subforest of \bar{x} from 2 to 4 is exactly the subtree rooted at 2. In general, subforests which correspond to subtrees are important special cases, which we can characterize in terms of *outermost right leafs*.

Definition 12 (outermost right leaf). Let X be a forest of size m = |X|. Further, let $i \in \{1, ..., m\}$. We define the outermost right leaf of i as:

$$rl_X(i) := \begin{cases} i & \text{if } \bar{x}_i = x \\ rl_X(i') & \text{if } \bar{x}_i = x(\bar{x}_{i,1}, \dots, \bar{x}_{i,R}) \text{ and } \bar{x}_{i'} = \bar{x}_{i,R} \end{cases}$$
(10)

Algorithm 2 An algorithm to retrieve the subforest from i to j of a forest X.

```
function SUBFOREST (A forest X = \bar{x}_1, \dots, \bar{x}_R, a start index i, an end index j, and a current index
k
    Y \leftarrow \epsilon.
    for r = 1, \ldots, R do
         k \leftarrow k + 1.
         if k > j then
              return (Y, k).
         else if k \geq i then
              (Y', k) \leftarrow \text{SUBFOREST}(\bar{\varrho}(\bar{x}_r), i, j, k).
              x \leftarrow \nu(\bar{x}_r).
              Y \leftarrow Y \oplus x(Y').
          else
              (Y', k) \leftarrow \text{SUBFOREST}(\bar{\varrho}(\bar{x}_r), i, j, k).
         end if
    end for
    return (Y, k).
end function
```

Again, consider the tree $\bar{x}=\mathtt{a}(\mathtt{b}(\mathtt{c},\mathtt{d}),\mathtt{e})$ from Figure 2. For this tree, we have $rl_{\bar{x}}(1)=5$, $rl_{\bar{x}}(2)=4$, $rl_{\bar{x}}(2)=4$, $rl_{\bar{x}}(3)=3$, $rl_{\bar{x}}(4)=4$, $rl_{\bar{x}}(5)=5$.

For the outermost right leafs it holds:

Theorem 6. Let X be a forest of size m = |X|. For any $i \in \{1, ..., m\}$ it holds:

$$rl_X(i) = i + |\bar{x}_i| \tag{11}$$

$$X[i, rl_X(i)] = \bar{x}_i \tag{12}$$

Proof. First, note that the pre-order algorithm (see definition 4) visits parents before children and left children before right children. Therefore, the largest index within a subtree must be the outermost right leaf.

Second, note that the pre-order algorithm visits all nodes in a subtree before leaving the respective subtree. Therefore, the outermost right leaf $rl_X(i)$ must be the index i plus the size of the subtree rooted at i, which proves the first claim.

The second claim follows because the subforest Algorithm 2 visits nodes in the same order as the pre-order algorithm and therefore $X[i, rl_X(i)] = X[i, i + |\bar{x}_i|] = \bar{x}_i$.

Now, we can define the edit distance between partial trees, which we call the *subforest edit distance*:

Definition 13 (Subforest edit distance). Let X and Y be forests over some alphabet \mathcal{X} , let c be a cost function over \mathcal{X} , let \bar{x}_k be an ancestor of \bar{x}_i in X, and let \bar{y}_l be an ancestor of \bar{y}_j in Y. Then, we define the subforest edit distance between the subforests $X[i, rl_X(k)]$ and $Y[j, rl_Y(l)]$ as

$$D_{c}(X[i, rl_{X}(k)], Y[j, rl_{Y}(l)]) := \min_{\bar{\delta} \in \Delta_{X}^{*}} \left\{ c(\bar{\delta}, X[i, rl_{X}(k)]) \middle| \bar{\delta}(X[i, rl_{X}(k)]) = Y[j, rl_{Y}(l)] \right\}$$
(13)

It directly follows that:

Theorem 7. Let X and Y be trees over some alphabet \mathcal{X} of size m = |X| and n = |Y| respectively. For every $i \in \{1, ..., m\}$ and $j \in \{1, ..., n\}$ we have:

$$D(X[i, rl_X(i)], Y[j, rl_Y(j)]) = d_c(\bar{x}_i, \bar{y}_j)$$
(14)

Proof. From Theorem 6 we know that $X[i, rl_X(i)] = \bar{x}_i$ and $Y[j, rl_Y(j)] = \bar{y}_j$. Therefore, we have

$$D(X[i,rl_X(i)],Y[j,rl_Y(j)]) := \min_{\bar{\delta} \in \Delta_{\mathcal{X}}^*} \left\{ c(\bar{\delta},\bar{x}_i) \middle| \bar{\delta}(\bar{x}_i) = \bar{y}_j \right\}$$

which corresponds exactly to the definition of $d_c(\bar{x}_i, \bar{y}_i)$.

Finally, we can go on to prove the arguably most important theorem for the TED, namely the recursive decomposition of the subforest edit distance:

Theorem 8. Let X and Y be non-empty forests over some alphabet X, let c be a cost function over X that is non-negative, self-equal, and conforms to the triangular inequality, let \bar{x}_k be an ancestor of \bar{x}_i in X, and let \bar{y}_l be an ancestor of \bar{y}_j in Y. Then it holds:

$$D_{c}(X[i,rl_{X}(k)],Y[j,rl_{Y}(l)]) = \min \left\{ c(x_{i},-) + D_{c}(X[i+1,rl_{X}(k)],Y[j,rl_{Y}(l)]), \\ c(-,y_{j}) + D_{c}(X[i,rl_{X}(k)],Y[j+1,rl_{Y}(l)]), \\ d_{c}(\bar{x}_{i},\bar{y}_{j}) + D_{c}(X[rl_{X}(i)+1,rl_{X}(k)],Y[rl_{Y}(j)+1,rl_{Y}(l)]) \right\}$$

$$(15)$$

Further it holds:

$$d_{c}(\bar{x}_{i}, \bar{y}_{j}) = \min\{c(x_{i}, -) + D_{c}(X[i + 1, rl_{X}(i)], Y[j, rl_{Y}(j)]),$$

$$c(-, y_{j}) + D_{c}(X[i, rl_{X}(i)], Y[j + 1, rl_{Y}(j)]),$$

$$c(x_{i}, y_{j}) + D_{c}(X[i + 1, rl_{X}(i)], Y[j + 1, rl_{Y}(j)])\}$$

$$(16)$$

Proof. We first show that an intermediate decomposition holds. In particular, we show that:

$$D_{c}(X[i,rl_{X}(k)],Y[j,rl_{Y}(l)]) = \min \left\{ c(x_{i},-) + D_{c}(X[i+1,rl_{X}(k)],Y[j,rl_{Y}(l)]), \\ c(-,y_{j}) + D_{c}(X[i,rl_{X}(k)],Y[j+1,rl_{Y}(l)]), \\ c(x_{i},y_{j}) + D_{c}(X[i+1,rl_{X}(i)],Y[j+1,rl_{Y}(j)]) + D_{c}(X[rl_{X}(i)+1,rl_{X}(k)],Y[rl_{Y}(j)+1,rl_{Y}(l)]) \right\}$$

$$(17)$$

Now, because we require that c is non-negative, self-equal, and conforms to the triangular inequality, we know that Theorem 4 holds, that is, we know that we can replace the cost of a cheapest edit script with the cost of a cheapest mapping. Let M be a cheapest mapping between the subtrees $X[i, rl_X(k)]$ and $Y[j, rl_Y(l)]$. Regarding i and j, only the following cases can occur:

- 1. i is not part of the mapping. In that case, x_i is deleted and we have $D_c(X[i, rl_X(k)], Y[j, rl_Y(l)])$ = $c(x_i, -) + D_c(X[i+1, rl_X(k)], Y[j, rl_Y(l)])$.
- 2. j is not part of the mapping. In that case, y_j is inserted and we have $c(-, y_j) + D_c(X[i, rl_X(k)], Y[j+1, rl_Y(l)])$.
- 3. Both i and j are part of the mapping. Let j' be the index i is mapped to and let i' be the index that is mapped to j, that is, $(i,j') \in M$ and $(i',j) \in M$. Because of the third constraint on mappings we know that $i \geq i' \iff j' \geq j$ and $i \leq i' \iff j' \leq j$. Now, consider the case that i' > i. In that case we know that j' < j. However, in that case, j' is not part of the subforest $Y[j, rl_Y(l)]$, because j is per definition the smallest index within the subforest. Therefore, (i, j') can not be part of a cheapest mapping between our two considered subforests.

Conversely, consider the case i' < i. This is also not possible because in that case i' is not part of the subforest $X[i, rl_X(k)]$. Therefore, it must hold that i' = i. However, this implies by the first constraint on mappings that j' = j. Therefore, $(i, j) \in M$. So we know that x_i is replaced with y_j , which implies that $D_c(X[i, rl_X(k)], Y[j, rl_Y(l)]) = c(x_i, y_j) + D_c(X[i+1, rl_X(k)], Y[j+1, rl_Y(l)])$.

However, if $(i, j) \in M$, the fourth constraint on mappings implies that all descendants of \bar{x}_i are mapped to descendants of \bar{y}_j . More specifically, for any $(i', j') \in M$ where \bar{x}_i is an ancestor of $\bar{x}_{i'}$ it must hold that \bar{y}_j is also an ancestor of $\bar{y}_{j'}$. Therefore, the subforest edit distance further decomposes into: $D_c(X[i, rl_X(k)], Y[j, rl_Y(l)]) = c(x_i, y_j) + D_c(X[i+1, rl_X(i)], Y[j+1, rl_Y(j)]) + D_c(X[rl_X(i)+1, rl_X(k)], Y[rl_Y(j)+1, rl_Y(l)])$.

Because we required that M is a cheapest mapping, the minimum of these three options must be the case, which yields Equation 17.

Using this intermediate result, we can now go on to prove Equations 15 and 16. In particular, Equation 16 holds because we find that:

$$\begin{split} d_c(\bar{x}_i,\bar{y}_j) = & D_c(X[i,rl_X(i)],Y[j,rl_Y(j)]) \\ = & \min \Big\{ c(x_i,-) + D_c(X[i+1,rl_X(i)],Y[j,rl_Y(j)]), \\ c(-,y_j) + D_c(X[i,rl_X(i)],Y[j+1,rl_Y(j)]), \\ c(x_i,y_j) + D_c(X[i+1,rl_X(i)],Y[j+1,rl_Y(j)]) + \\ D_c(X[rl_X(i)+1,rl_X(i)],Y[rl_Y(j)+1,rl_Y(j)]) \Big\} \end{split}$$

where $D_c(X[rl_X(i)+1,rl_X(i)],Y[rl_Y(j)+1,rl_Y(j)])=0$ because the two input forests are empty.

With respect to Equation 15, we observe that one way to transform the subforest $X[i, rl_X(k)]$ into the subforest $Y[j, rl_Y(l)]$ is to transform the subforest $X[i, rl_X(i)]$ into the subforest $Y[j, rl_Y(j)]$, and then the subforest $X[rl_X(i) + 1, rl_X(k)]$ into the subforest $Y[rl_Y(j) + 1, rl_Y(l)]$. Therefore, we obtain:

$$D_{c}(X[i, rl_{X}(k)], Y[j, rl_{Y}(l)])$$

$$\leq D_{c}(X[i, rl_{X}(i)], Y[j, rl_{Y}(j)]) + D_{c}(X[rl_{X}(i) + 1, rl_{X}(k)], Y[rl_{Y}(j) + 1, rl_{Y}(l)])$$

$$= d_{c}(\bar{x}_{i}, \bar{y}_{i}) + D_{c}(X[rl_{X}(i) + 1, rl_{X}(k)], Y[rl_{Y}(j) + 1, rl_{Y}(l)]),$$

$$(18)$$

Further, we observe that

$$d_c(\bar{x}_i, \bar{y}_i) \le c(x_i, y_i) + D_c(X[i+1, rl_X(i)], Y[j+1, rl_Y(j)]), \tag{19}$$

because this is only one of the three cases in Equation 16.

Now, note that the first two cases in Equations 15 and 17 are the same. Finally, consider that the last case of Equation 17. In that case, we can conclude that:

$$\begin{split} &d_c(\bar{x}_i,\bar{y}_j) + D_c(X[rl_X(i)+1,rl_X(k)],Y[rl_Y(j)+1,rl_Y(l)]) \\ &19 \\ &\leq c(x_i,y_j) + D_c(X[i+1,rl_X(i)],Y[j+1,rl_Y(j)]) + \\ &D_c(X[rl_X(i)+1,rl_X(k)],Y[rl_Y(j)+1,rl_Y(l)]) \\ &\frac{17}{2}D_c(X[i,rl_X(k)],Y[j,rl_Y(l)]) \\ &18 \\ &\leq d_c(\bar{x}_i,\bar{y}_j) + D_c(X[rl_X(i)+1,rl_X(k)],Y[rl_Y(j)+1,rl_Y(l)]) \end{split}$$

which implies that Equation 15 holds.

As an example, consider the trees $\bar{x} = \mathsf{a}(\mathsf{b}(\mathsf{c},\mathsf{d}),\mathsf{e})$ and $\bar{y} = \mathsf{f}(\mathsf{g})$ from Figure 2 and the subforest edit distance $D_c(\bar{x}[2,5],\bar{y}[1,2]) = D_c((\mathsf{b}(\mathsf{c},\mathsf{d}),\mathsf{e}),\mathsf{f}(\mathsf{g}))$. According to Equation 15, we can decompose this in three ways, corresponding to the options to delete b , replace b with f , and insert f . In particular, we obtain the distance $c(\mathsf{b},-) + D_c((\mathsf{c},\mathsf{d},\mathsf{e}),\mathsf{f}(\mathsf{g}))$ for the deletion, $d_c(\mathsf{b}(\mathsf{c},\mathsf{d}),\mathsf{f}(\mathsf{g})) + D_c(\mathsf{e},\mathsf{e})$ for the replacement, and $c(-,\mathsf{f}) + D_c((\mathsf{b}(\mathsf{c},\mathsf{d}),\mathsf{e}),\mathsf{g})$ for the insertion (also refer to Figure 5).

Now, consider the replacement option. According to Equation 16, we can decompose the TED $d_c(b(c,d), f(g))$ in three ways, corresponding to the options to delete b, replace b with f, and insert f. In particular, we obtain the distance $c(b,-) + D_c((d,e), f(g))$ for the deletion, $c(b,f) + D_c((d,e),g)$ for the replacement, and $c(-,f) + D_c(b(c,d),g)$ for the insertion (also refer to Figure 5).

For an efficient algorithm for the TED, we are missing only one last ingredient, namely a valid base case for empty forests. This is easy enough to obtain:

Theorem 9. Let X and Y be forests over some alphabet \mathcal{X} , let c be a cost function \mathcal{X} that is non-negative, self-equal, and conforms to the triangular inequality, let \bar{x}_k be an ancestor of \bar{x}_i in X, and let \bar{y}_l be an ancestor of \bar{y}_j in Y. Then it holds:

$$D_c(\epsilon, \epsilon) = 0 \tag{20}$$

$$D_c(X[i, rl_X(k)], \epsilon) = c(x_i, -) + D_c(X[i+1, rl_X(k)], \epsilon)$$
(21)

$$D_c(\epsilon, Y[j, rl_Y(l)]) = c(-, y_j) + D_c(\epsilon, Y[j+1, rl_Y(l)])$$
(22)

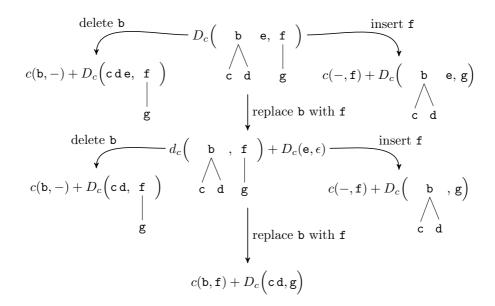


Figure 5: An illustration of the decompositions in Equation 15 and 16.

Proof. Because c non-negative and self-equal, the cheapest script to transform an empty forest into an empty forest is to do nothing. Further, because c conforms to the trinagular inequality, the cheapest script to transform a non-empty forest into an empty forest is to delete all nodes. Finally, the cheapest script to transform an empty forest into a non-empty one is to insert all nodes.

Now, we are able to construct an efficient algorithm for the TED. We just need to iterate over all possible pairs of subtrees in both input trees and compute the TED between these subtrees. For this, we require the subforest edit distance for all pairs of subforests in these subtrees. We store intermediate results for the subforest edit distance in an array D and intermediate results for the subtree edit distance in an array d. Finally, the edit distance between the whole input trees will be in the first entry of d.

Theorem 10. The output of Algorithm 3 is the TED. Further, Algorithm 3 runs in $\mathcal{O}(m^2 \cdot n^2)$ time and in $\mathcal{O}(m \cdot n)$ space complexity where m and n are the sizes of the input trees.

Proof. Each computational step is justified by one of the equations proven before (refer to the comments in the pseudo-code). Therefore, the output of the algorithm is correct. Finally, the algorithm runs in $\mathcal{O}(m^2 \cdot n^2)$ because two of the nested for-loops run at most m times and two of the loops run at most n times. Note that this bound is tight because the worst case does occur for the trees shown in Figure 6. Regarding space complexity, we note that we maintain two matrices, d and D, each with $\mathcal{O}(m \cdot n)$ entries.

As Zhang and Shasha (1989) point out, we can be even more efficient in our algorithm if we re-use already computed subforest edit distances. In particular, we can re-use the subforest edit distance whenever the outermost right leaf is equal:

Theorem 11. Let X and Y be forests over some alphabet \mathcal{X} , let c be a cost function over \mathcal{X} , let \bar{x}_k be an ancestor of \bar{x}_i in X such that $rl_X(k) = rl_X(i)$, and let \bar{y}_l be an ancestor of \bar{y}_j such that $rl_Y(l) = rl_Y(j)$. Then it holds for all i such that \bar{x}_k is an ancestor of \bar{x}_i and all j such that \bar{y}_l is an ancestor of \bar{y}_j :

$$D_c(X[i, rl_X(k)], Y[j, rl_Y(l)]) = D_c(X[i, rl_X(k')], Y[j, rl_Y(l')])$$
(23)

Proof. Because we required that $rl_X(k) = rl_X(i)$ and $rl_Y(l) = rl_Y(j)$, this follows directly.

Algorithm 3 An efficient algorithm for the TED. Note that this algorithm is not yet the most efficient one, but a proto-version of the actual TED algorithm of Zhang and Shasha (1989) which is shown later as Algorithm 5. The algorithm iterates over all subtrees of \bar{x} and \bar{y} and computes the tree edit distance for them based on the forest edit distances between all subforests of the respective subtrees.

```
function TREE-EDIT-DISTANCE(Two input trees \bar{x} and \bar{y}, a cost function c.)
     m \leftarrow |\bar{x}|, n \leftarrow |\bar{y}|.
     \boldsymbol{d} \leftarrow m \times n matrix of zeros.
                                                                                                                                       \triangleright d_{i,j} = d_c(\bar{x}_i, \bar{y}_j).
     \mathbf{D} \leftarrow (m+1) \times (n+1) matrix of zeros.
                                                                                                       \triangleright \mathbf{D}_{i,j} = D_c(X[i, rl_{\bar{x}}(k)], Y[j, rl_{\bar{y}}(l)]).
     for k \leftarrow m, \dots, 1 do
           for l \leftarrow n, \dots, 1 do
                                                                                                                                              ⊳ Equation 20
                 \boldsymbol{D}_{rl_X(k)+1,rl_Y(l)+1} \leftarrow 0.
                 for i \leftarrow rl_X(k), \ldots, k do
                       D_{i,rl_Y(l)+1} \leftarrow D_{i+1,rl_Y(l)+1} + c(x_i, -).
                                                                                                                                              ⊳ Equation 21
                 end for
                 for j \leftarrow rl_Y(l), \ldots, l do
                      D_{rl_X(k)+1,j} \leftarrow D_{rl_X(k)+1,j+1} + c(-,y_j).
                                                                                                                                              ▶ Equation 22
                 end for
                 for i \leftarrow rl_X(k), \ldots, k do
                       for j \leftarrow rl_Y(l), \ldots, l do
                            if rl_{\bar{x}}(i) = rl_{\bar{x}}(k) \wedge rl_{\bar{y}}(j) = rl_{\bar{y}}(l) then
                                  \mathbf{D}_{i,j} \leftarrow \min{\{\mathbf{D}_{i+1,j} + c(x_i, -), \}}
                                      \hat{D}_{i,j+1} + \hat{c}(-,y_i),
                                      D_{i+1,j+1} + c(x_i, y_j).
                                                                                                                                              ⊳ Equation 16
                                  oldsymbol{d}_{i,j} \leftarrow 	ilde{oldsymbol{D}}_{i,j}.
                            else
                                  \boldsymbol{D}_{i,j} \leftarrow \min\{\boldsymbol{D}_{i+1,j} + c(x_i, -), \\ \boldsymbol{D}_{i,j+1} + c(-, y_j),
                                      D_{rl_{\bar{x}}(i)+1,rl_{\bar{y}}(j)+1}+d_{i,j}.
                                                                                                                                              ▷ Equation 15
                            end if
                      end for
                 end for
           end for
     end for
     return d_{1,1}.
end function
```

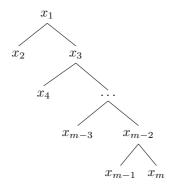


Figure 6: An example tree structure which yields the worst-case runtime of $\mathcal{O}(m^2 \cdot n^2)$ in Algorithm 3 as well as Algorithm 5.

Therefore, we can make our algorithm faster by letting the two outer loops only run over nodes for which we can *not* re-use the subforest edit distance. Those nodes are the so-called *keyroots* of our input trees.

Definition 14 (keyroots). Let X be a forest over some alphabet \mathcal{X} and let \bar{x}_i be a leaf in X. We define the *keyroot* of \bar{x}_i as

$$k_X(i) = \min\{k | rl_X(k) = i\}$$
(24)

We define the keyroots of X, denoted as $\mathcal{K}(X)$ as the set of keyroots for all leaves of X.

For example, if we inspect tree $\bar{x} = \mathsf{a}(\mathsf{b}(\mathsf{c},\mathsf{d}),\mathsf{e})$ from Figure 2 our leaves are $\bar{x}_3 = \mathsf{c}$, $\bar{x}_4 = \mathsf{d}$, and $\bar{x}_5 = \mathsf{e}$. The corresponding key roots are $k_{\bar{x}}(3) = 3$, $k_{\bar{x}}(4) = 2$, and $k_{\bar{x}}(5) = 1$. Accordingly, the set of keyroots $\mathcal{K}(\bar{x})$ is $\{1, 2, 3\}$.

Computing the keyroots is possible using Algorithm 4.

Algorithm 4 An algorithm to compute the key roots of a forest.

```
\begin{aligned} & \textbf{function} \text{ KEYROOTS}(\text{A forest } X.) \\ & m \leftarrow |X|. \\ & R \leftarrow \emptyset. \\ & \mathcal{K} \leftarrow \epsilon. \\ & \textbf{for } i \leftarrow 1, \ldots, m \textbf{ do} \\ & rl \leftarrow rl_X(i). \\ & \textbf{if } rl \notin R \textbf{ then} \\ & R \leftarrow R \cup \{rl\}. \\ & \mathcal{K} \leftarrow \mathcal{K} \oplus i. \\ & \textbf{end if} \\ & \textbf{end for} \\ & \textbf{return } \mathcal{K}. \\ & \textbf{end function} \end{aligned}
```

This yields the TED Algorithm 5 of Zhang and Shasha (1989).

Theorem 12. Algorithm 5 computes the TED. Further, Algorithm 5 runs in $\mathcal{O}(m^2 \cdot n^2)$ and has $\mathcal{O}(m \cdot n)$ space complexity.

Proof. The runtime proof is simple: Because we use a subset of outer loop iterations compared to Algorithm 3, we are at most as slow. Still, the runtime bound is tight, because the set of keyroots is per definition as large as the set of leaves of a tree, and the number of leaves of a tree can grow linearly with the size of a tree, as is the case in Figure 6. The space requirements are the same as for Algorithm 3.

Further, Algorithm 5 still computes the same result as Algorithm 3, because according to Theorem 11 the same subforest edit distances are computed as before.

As an example, consider the trees $\bar{x}=a(b(c,d),e)$ and $\bar{y}=f(g)$ from Figure 2. The TED algorithm first considers the edit distance between the subtrees $\bar{x}_3=c$ and $\bar{y}_2=f(g)$, which can be computed based on the subforest edit distances $D_c(\bar{x}[4,3],\bar{y}[3,2])=D_c(\epsilon,\epsilon)=0$, $D_c(\bar{x}[4,3],\bar{y}[2,2])=D_c(\epsilon,g)=1$, $D_c(\bar{x}[4,3],\bar{y}[1,2])=D_c(\epsilon,f(g))=2$, $D_c(\bar{x}[3,3],\bar{y}[3,2])=D_c(c,\epsilon)=1$, and $D_c(\bar{x}[3,3],\bar{y}[2,2])=D_c(c,g)=1$. Based on these intermediate results, we can infer that $d_c(\bar{x}_3,\bar{y}_1)=D_c(\bar{x}[3,3],\bar{y}[1,2])=2$. Note that this calculation also yields the edit distance between the subtrees $\bar{x}_3=c$ and $\bar{y}_2=g$ as an intermediate result (also refer to Figure 7, middle). Next, we compute the subtree edit distance between $\bar{x}_2=b(c,d)$ and $\bar{y}_1=f(g)$, which also yields the subtree edit distances $d_c(\bar{x}_4,\bar{y}_2),d_c(\bar{x}_4,\bar{y}_1),$ and $d_c(\bar{x}_2,\bar{y}_2)$ as intermediate results (see Figure 7, middle). Finally, we can compute the subtree edit distance between $\bar{x}_1=\bar{x}$ and $\bar{y}_1=\bar{y}$ (see Figure 7, bottom), which turns out to be 5.

This concludes our description of the edit distance itself. However, in many situations it is not only interesting to know the size of the edit distance, but also which mapping (and which edit script) corresponds to the edit distance. This is the topic of *backtracing*.

Algorithm 5 The $\mathcal{O}(m^2 \cdot n^2)$ TED algorithm of Zhang and Shasha (1989). The algorithm iterates over all subtrees of \bar{x} and \bar{y} rooted at key roots and computes the TED for them based on the forest edit distances between all subforests of the respective subtrees. Refer to our project web site for a reference implementation.

```
function TREE-EDIT-DISTANCE (Two input trees \bar{x} and \bar{y}, a cost function c.)
     m \leftarrow |\bar{x}|, n \leftarrow |\bar{y}|.
     \mathcal{K}(\bar{x}) \leftarrow \text{KEYROOTS}(\bar{x}).
     \mathcal{K}(\bar{y}) \leftarrow \text{KEYROOTS}(\bar{y}).
                                                                                                                                     \triangleright d_{i,j} = d_c(\bar{x}_i, \bar{y}_j).
     d \leftarrow m \times n matrix of zeros.
     D \leftarrow (m+1) \times (n+1) matrix of zeros.
                                                                                                      \triangleright D_{i,j} = D_c(X[i, rl_{\bar{x}}(k)], Y[j, rl_{\bar{y}}(l)]).
     for k \in \mathcal{K}(\bar{x}) in descending order do
           for l \in \mathcal{K}(\bar{y}) in descending order do
                 \boldsymbol{D}_{rl_X(k)+1,rl_Y(l)+1} \leftarrow 0.
                                                                                                                                             ⊳ Equation 20
                 for i \leftarrow rl_X(k), \ldots, k do
                      D_{i,rl_Y(l)+1} \leftarrow D_{i+1,rl_Y(l)+1} + c(x_i, -).
                                                                                                                                             ⊳ Equation 21
                 end for
                for j \leftarrow rl_Y(l), \ldots, l do
                      D_{rl_X(k)+1,j} \leftarrow D_{rl_X(k)+1,j+1} + c(-,y_j).
                                                                                                                                             ⊳ Equation 22
                end for
                for i \leftarrow rl_X(k), \ldots, k do
                      for j \leftarrow rl_Y(l), \ldots, l do
                            if rl_{\bar{x}}(i) = rl_{\bar{x}}(k) \wedge rl_{\bar{y}}(j) = rl_{\bar{y}}(l) then
                                  D_{i,j} \leftarrow \min\{D_{i+1,j} + c(x_i, -),
                                 D_{i,j+1} + c(-, y_j),

D_{i+1,j+1} + c(x_i, y_j).

d_{i,j} \leftarrow D_{i,j}.
                                                                                                                                             ▶ Equation 16
                            else
                                  \boldsymbol{D}_{i,j} \leftarrow \min\{\boldsymbol{D}_{i+1,j} + c(x_i, -), \\ \boldsymbol{D}_{i,j+1} + c(-, y_j),
                                      D_{rl_{\bar{x}}(i)+1,rl_{\bar{y}}(j)+1}+d_{i,j}.
                                                                                                                                             ▶ Equation 15
                            end if
                      end for
                end for
           end for
     end for
     return d_{1,1}.
end function
```

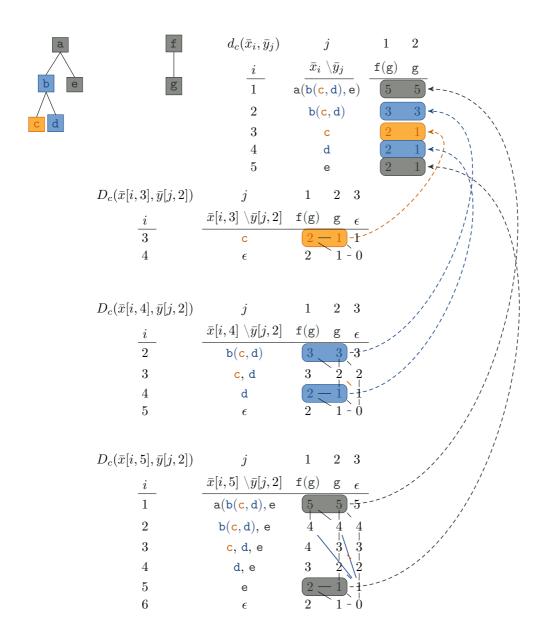


Figure 7: An illustration of the TED Algorithm 5 of Zhang and Shasha (1989) for the two input trees from Figure 2. Nodes with the same right outermost leaf are shown in the same color. For these nodes, the subforest edit distance is re-used. Top right: The TED between all subtrees of the input trees. Bottom: The subforest edit distances for all key root pairs. All entries which correspond to a subtree edit distance are highlighted in color and linked with dashed arrows to the corresponding entries in the subtree edit distance table at the top right. Co-optimal mappings are indicated by solid lines linking the entries of the dynamic programming table to the distances they are decomposed into.

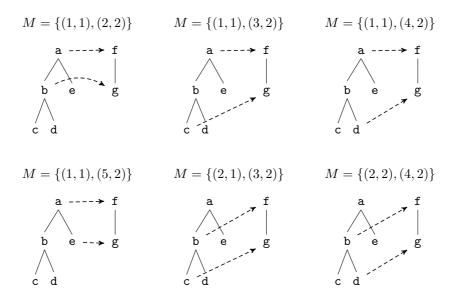


Figure 8: All six co-optimal mappings between the trees $\bar{x} = a(b(c,d),e)$ and $\bar{y} = f(g)$ from Figure 7.

4 Backtracing and Co-Optimal Mappings

Now that we have computed the TED, the next question is: Which edit script corresponds to the TED? We have answered this question in part in Algorithm 1, which transforms a mapping into an edit script with the same cost. Therefore, we can re-phrase the question: Which *mapping* corresponds to the TED? Zhang and Shasha (1989) only hint at an answer in their own paper. Here, we shall analyze this question in detail. We start by phrasing more precisely what we are looking for:

Definition 15 (Co-Optimal Mappings). Let X and Y be forests over some alphabet \mathcal{X} and let c be a cost function over \mathcal{X} . We define a co-optimal mapping M as a tree mapping between X and Y such that $c(M,X,Y) = \min_{M'} c(M',X,Y)$.

For example, all co-optimal mappings for the trees in Figure 7 are shown in Figure 8. Unfortunately, listing all co-optimal mappings is infeasible in general, as the following theorem demonstrates:

Theorem 13. Let $\bar{\mathbf{a}}(1)$ be the tree \mathbf{a} over the alphabet $\mathcal{X} = \{\mathbf{a}\}$ and let $\bar{\mathbf{a}}(m) = \mathbf{a}(\bar{\mathbf{a}}(m-1))$. Then, for any metric cost function c over \mathcal{X} and any $m \in \mathbb{N}$ which is divisible by 2 it holds: There are $m!/[(m/2)!]^2$ co-optimal mappings between the trees $\bar{x} = \bar{\mathbf{a}}(m)$ and $\bar{y} = \bar{\mathbf{a}}(m/2)$. Further, this number is larger than $\frac{\sqrt{2\pi}}{e^2} \cdot \frac{2^{m+1}}{\sqrt{m}}$.

Proof. Because c is metric, we have $c(\mathtt{a},\mathtt{a})=0$ and $c(\mathtt{a},-)>0$. Therefore, we want to replace as many \mathtt{a} with \mathtt{a} as possible to reduce the cost. At most, we can replace m/2 \mathtt{a} with \mathtt{a} , because $|\bar{\mathtt{a}}(m/2)|=m/2$. Therefore, this corresponds to choosing m/2 nodes from \bar{x} which are mapped to the m/2 nodes from \bar{y} . As we know from combinatorics, there are

$$\binom{m}{\frac{m}{2}} = \frac{m!}{(\frac{m}{2}!)^2}$$

ways to choose m/2 from m elements. Using Stirling's approximation, we then obtain the following lower bound:

$$\frac{m!}{(\frac{m!}{2}!)^2} \geq \frac{\sqrt{2\pi} \cdot m^{m+\frac{1}{2}} \cdot e^{-m}}{(e \cdot (\frac{m}{2})^{\frac{m}{2}+\frac{1}{2}} \cdot e^{-\frac{m}{2}})^2} = \frac{\sqrt{2\pi}}{e^2} \cdot \frac{m^{m+\frac{1}{2}}}{\frac{m}{2})^{m+\frac{1}{2}} \cdot \sqrt{\frac{m}{2}}} \cdot \frac{e^{-m}}{e^{-m}} = \frac{\sqrt{2\pi}}{e^2} \cdot \frac{2^{m+\frac{1}{2}}}{\sqrt{\frac{m}{2}}} = \frac{\sqrt{2\pi}}{e^2} \cdot \frac{2^{m+1}}{\sqrt{m}}$$

While it is therefore infeasible to list *all* co-optimal mappings, it is still possible to return *some* of the co-optimal mappings. In particular, it turns out that constructing a co-optimal mapping corresponds to finding a path in a graph which we call the *co-optimal edit graph*. First, we define a general graph as follows:

Definition 16 (Directed Acyclic Graph (DAG)). Let V be some set, and let $E \subseteq V \times V$. Then we call $\mathcal{G} = (V, E)$ a graph, V the nodes of \mathcal{G} and E the edges of \mathcal{G} . We call \mathcal{G} a directed acyclic graph (DAG) if there exists a total ordering relation < on V, such that for all edges $(u, v) \in E$ it holds: u < v, i.e. edges occur only from lower nodes to higher nodes in the ordering.

We then define our co-optimal edit graph as follows.

Definition 17 (Co-optimal Edit Graph). Let X and Y be forests over some alphabet \mathcal{X} and let c be a cost function over \mathcal{X} . Then, we define the *co-optimal edit graph* between X and Y according to c as the graph $\mathcal{G}_{c,X,Y} = (V,E)$ with nodes V and edges E as follows.

If $X = \epsilon$ and $Y = \epsilon$ we define $V := \{(1, 1, 1, 1)\}$ and $E := \emptyset$.

If $X = \epsilon$ but $Y \neq \epsilon$ we define $V := \{(1, 1, 1, j) | j \in \{1, \dots, |Y| + 1\} \}$ and $E := \{((1, 1, 1, j), (1, 1, 1, j + 1)) | j \in \{1, \dots, |Y| \} \}.$

If $X \neq \epsilon$ but $Y = \epsilon$ we define $V := \{(1, i, 1, 1) | i \in \{1, ..., |X| + 1\}\}$ and $E := \{((1, i, 1, 1), (1, i + 1, 1, 1)) | i \in \{1, ..., |X|\}\}$.

If neither forest is empty, we define:

$$V := \left\{ (k, i, l, j) \middle| k \in \mathcal{K}(X), i \in \{k, \dots, rl_X(k) + 1\}, l \in \mathcal{K}(Y), j \in \{l, \dots, rl_Y(l) + 1\} \right\}$$
(25)

$$E := \left\{ \left((k, i, l, j), (k, i + 1, l, j) \right) \middle| D_c \left(X[i, rl_X(k)], Y[j, rl_Y(l)] \right) \right\}$$

$$= c(x_i, -) + D_c(X[i+1, rl_X(k)], Y[j, rl_Y(l)]) \}$$
(26)

$$\Big\{ \Big((k,i,l,j), (k,i,l,j+1) \Big) \Big| D_c \Big(X[i,rl_X(k)], Y[j,rl_Y(l)] \Big)$$

$$= c(-, y_j) + D_c(X[i, rl_X(k)], Y[j+1, rl_Y(l)]) \} \cup$$
(27)

$$\Big\{ \Big((k,i,l,j), (k,i+1,l,j+1) \Big) \Big| D_c \Big(X[i,rl_X(k)], Y[j,rl_Y(l)] \Big)$$

$$= c(x_i, y_i) + D_c(X[i+1, rl_X(k)], Y[j+1, rl_Y(l)])$$

$$\wedge rl_X(i) = rl_X(k) \wedge rl_Y(j) = rl_Y(l) \} \cup \tag{28}$$

$$\Big\{ \Big((k,i,l,j), (k,i+1,l,j+1) \Big) \Big| D_c \Big(X[i,rl_X(k)], Y[j,rl_Y(l)] \Big)$$

$$= c(x_i, y_i) + D_c(X[i+1, rl_X(k)], Y[j+1, rl_Y(l)])$$

$$\wedge c(x_i, y_j) = c(x_i, -) + c(-, y_j) \Big\} \cup$$
 (29)

$$\left\{ ((k,i,l,j), (k_X(i), i+1, k_Y(j), j+1)) \middle| D_c(X[i,rl_X(k)], Y[j,rl_Y(l)]) \right\}$$

$$= D_c(\bar{x}^i, \bar{y}^j) + D_c(X[rl_X(i) + 1, rl_X(k)], Y[rl_Y(j) + 1, rl_Y(l)])$$

$$\wedge \left(rl_X(i) \neq rl_X(k) \vee rl_Y(j) \neq rl_Y(l) \right) \wedge c(x_i, y_j) < c(x_i, -) + c(-, y_j) \right\} \cup \tag{30}$$

$$\Big\{ \big((k, rl_X(k) + 1, l, rl_Y(l) + 1), (\mathbf{k}_X(rl_X(k) + 1), rl_X(k) + 1, \mathbf{k}_Y(rl_Y(l) + 1), rl_Y(l) + 1) \big) \Big|$$

$$rl_X(k) + 1 \le |X| \land rl_Y(l) + 1 \le |Y|$$
 (31)

$$\left\{ \left((k, rl_X(k) + 1, l, |Y| + 1), (k_X(rl_X(k) + 1), rl_X(k) + 1, 1, |Y| + 1) \right) \middle| rl_X(k) + 1 \le |X| \right\} \cup (32)$$

$$\left\{ \left((k, |X| + 1, l, rl_Y(l) + 1), (1, |X| + 1, k_Y(rl_Y(l) + 1), rl_Y(l) + 1) \right) \middle| rl_Y(l) + 1 \le |Y| \right\}$$
(33)

As this definition is quite extensive, we shall explain it in a bit more detail. The nodes of the cooptimal edit graph are, essentially, the entries of the dynamic programming matrix D of Algorithm 5. Given that this matrix needs to be computed for every combination of keyroots $(k, l) \in \mathcal{K}(X) \times \mathcal{K}(Y)$,

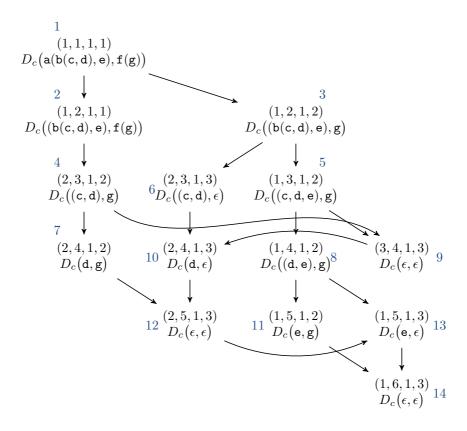


Figure 9: An excerpt of the co-optimal edit graph between the trees \bar{x} and \bar{y} from Figure 2. The figure only shows nodes which are reachable from (1,1,1,1). Further, to support clarity, the nodes are labelled with indices and with the corresponding subforest edit distance. The indices in blue mark the order according to the ordering relationship < as defined in Theorem 14

we need four indices to identify a position in the dynamic programming matrix D uniquely, namely the keyroot indices and the matrix indices, leading to a quartuple (k, i, l, j). Now, with respect to the edges, Equation 26 defines the edges corresponding to deletions (the first case in Equations 15 and 16), Equation 27 defines the edges corresponding to insertions (the second case in Equations 15 and 16), Equation 28 defines the edges corresponding to replacements within a subtree (the third case in Equation 16), and Equation 30 defines the edges corresponding to replacements of entire subtrees (the third case in Equation 15).

The remaining edges cover special cases. In particular, Equation 29 covers the case where all options, deletion, insertion, and replacement are co-optimal, and Equations 31, 32, and 33 cover cases where we are at the end of the dynamic programming matrix for a subtree and need to continue the computation in the dynamic programming matrix for a larger subtree which includes the current subtree.

As an example, consider the co-optimal edit graph between the trees $\bar{x} = a(b(c, d), e)$ and $\bar{y} = f(g)$ from Figure 2. An excerpt of this graph is shown in Figure 9.

An important insight regarding the co-optimal edit graph is that it is acyclic.

Theorem 14. Let X and Y be forests over some alphabet \mathcal{X} , let c be a cost function over \mathcal{X} , and let $\mathcal{G}_{c,X,Y}$ be the co-optimal edit graph with respect to X, Y, and c. Then, $\mathcal{G}_{c,X,Y}$ is a directed acyclic graph with the ordering relation (k,i,l,j) < (k',i',l',j') if and only if i < i', or i = i' and j < j', or i = i' and j = j' and k > k', or i = i' and j = j' and k = k' and k = k' and k = k' and k = k'.

Proof. Follows directly from the definition of the edges.

The ordering for the example graph in Figure 9 is displayed as blue indices.

Each edge in the co-optimal edit graph corresponds to an edit which could be used in a co-optimal edit script. Accordingly, we should be able to join edges in the co-optimal edit graph together, such that we obtain a complete, co-optimal edit script. This notion of joining edges is captured by the notion of a *path*.

Definition 18 (path). Let $\mathcal{G} = (V, E)$ be a graph. A path p from $u \in V$ to $v \in V$ is defined as a sequence of nodes $p = v_0, \ldots, v_T$, such that $v_0 = u$, $v_T = v$, and for all $t \in \{1, \ldots, T\} : (v_{t-1}, v_t) \in E$. If a path from u to v exists, we call v reachable from u.

Note that our definition permits trivial paths of length T=0 from any node to itself. Next, we define the corresponding mapping to a path in the co-optimal edit graph:

Definition 19 (corresponding mapping). Let X and Y be forests over some alphabet \mathcal{X} , let c be a cost function over \mathcal{X} , and let $\mathcal{G}_{c,X,Y} = (V,E)$ be the co-optimal edit graph with respect to X,Y, and c. Further, let v_0, \ldots, v_T be a path from 1 to |V| in $\mathcal{G}_{c,X,Y}$. Then, we define the *corresponding mapping* $M_{X,Y}(v_0,\ldots,v_T)$ for path v_0,\ldots,v_T as $M=\emptyset$ if X or Y are empty and as follows otherwise:

$$M_{X,Y}(v_0, \dots, v_T) = \{(i, j) | v_t = (k, i, l, j), v_{t+1} = (k', i+1, l', j+1) \text{ for any } k, l, k', l', t\}$$
 (34)

Consider the example path p = (1, 1, 1, 1), (1, 2, 1, 2), (1, 3, 1, 2), (3, 4, 1, 3), (2, 4, 1, 3), (2, 5, 1, 3), (1, 5, 1, 3), (1, 6, 1, 3) in Figure 9. The corresponding mapping $M_{a(b(c,d),e),f(g)}(p)$ would be $\{(1,1),(3,2)\}$, corresponding to the replacement edges on the path.

Given these definitions, we can go on to show our key theorem for backtracing, namely that the corresponding mapping for every path through the co-optimal edit graph is a co-optimal mapping, and that every co-optimal mapping corresponds to a path through the co-optimal edit graph.

Theorem 15. Let X and Y be forests over some alphabet \mathcal{X} , let c be a cost function over \mathcal{X} that is non-negative, self-equal, and conforms to the trinagular inequality, and let $\mathcal{G}_{c,X,Y}$ be the co-optimal edit graph with respect to X, Y, and c. Then it holds:

- 1. For all paths p from (1,1,1,1) to (1,|X|+1,1,|Y|+1) in $\mathcal{G}_{c,X,Y}$, the corresponding mapping $M_{X,Y}(p)$ is a co-optimal mapping between X and Y.
- 2. For all co-optimal mappings M between X and Y, there exists at least one path p from (1, 1, 1, 1) to (1, |X| + 1, 1, |Y| + 1) in $\mathcal{G}_{c,X,Y}$ such that $M_{X,Y}(p) = M$.

Proof. As a shorthand, we will call a path from (1,1,1,1) to (1,|X|+1,1,|Y|+1) in a co-optimal edit graph $\mathcal{G}_{c,X,Y}$ a path through that graph.

We start by considering the trivial cases of empty forests. If $X = \epsilon$ or $Y = \epsilon$, the only co-optimal mapping is $M = \emptyset$. It remains to show that, in these cases, the co-optimal edit graph contains only paths which correspond to this mapping.

 $X = \epsilon$ and $Y = \epsilon$: In this case, we obtain $V = \{(1, 1, 1, 1)\}$ and $E = \emptyset$. Accordingly, the trivial path p = (1, 1, 1, 1) is the only possible path through $\mathcal{G}_{c,X,Y}$ and it does indeed hold $M_p = \emptyset$.

 $X = \epsilon$ and $Y \neq \epsilon$: In this case, we obtain $V = \{(1, 1, 1, j) | j \in \{1, \dots, |Y| + 1\}\}$ and

$$E = \{ ((1,1,1,j), (1,1,1,j+1)) | j \in \{1,\dots,|Y|\} \}.$$

Accordingly, the only possible path through $\mathcal{G}_{c,X,Y}$ is $p = (1,1,1,1), (1,1,1,2), \ldots, (1,1,1,|Y|+1)$. And indeed it holds $M_p = \emptyset$.

 $X \neq \epsilon$ and $Y = \epsilon$: In this case, we obtain $V = \{(1, i, 1, 1) | i \in \{1, \dots, |X| + 1\}\}$ and

$$E = \{ ((1, i, 1, 1), (1, i + 1, 1, 1)) | i \in \{1, \dots, |X|\} \}.$$

Accordingly, the only possible path through $\mathcal{G}_{c,X,Y}$ is $p = (1,1,1,1), (1,2,1,1), \dots, (1,|X| + 1,1,1)$. And indeed it holds $M_p = \emptyset$.

It remains to show both claims for the case of non-empty forests. For both claims, we apply an induction over the added size of both input forests, for which the base case is already provided by our considerations above. For the induction, we assume that both claims hold for inputs forests X and Y with $|X| + |Y| \le k$.

Now, we consider two input forests X and Y with m+n=k+1, where m=|X| and n=|Y|. Regarding the first claim, let $p=v_0,\ldots,v_T$ be a path through $\mathcal{G}_{c,X,Y}$, let X':=X[2,|X|], let Y':=Y[2,|Y|], and consider the following cases regarding v_1 .

- $v_1 = (1,2,1,1)$: In this case, it must hold $D_c(X,Y) = c(x_1,-) + D_c(X',Y)$, otherwise $(v_0,v_1) \notin E$. Now, if X' is empty, then p must have the form $p = (1,1,1,1), (1,2,1,1), (1,2,1,2), \ldots$, (1,2,1,|Y|+1), and \emptyset must be a co-optimal mapping between X' and Y. Accordingly, $\emptyset = M_p$ must also be a co-optimal mapping between X and Y, because $c(\emptyset,X,Y) = c(x_1,-) + c(\emptyset,X',Y) = c(x_1,-) + D_c(X',Y) = D_c(X,Y)$.
 - If X' is not empty, then the path v_1, \ldots, v_T is isomorphic to a path from (1,1,1,1) to (1,|X'|+1,1,|Y|+1) in $\mathcal{G}_{c,X',Y}$. Accordingly, per induction hypothesis, $M_{p'}$ is a co-optimal mapping between X' and Y. Further, we obtain per construction $M_p = \{(i+1,j)|(i,j) \in M_{p'}\}$. Accordingly, it holds: $c(M_p,X,Y) = c(x_1,-) + c(M_{p'},X',Y) = c(x_1,-) + D_c(X',Y) = D_c(X,Y)$, which means that M_p is co-optimal, as claimed.
- $v_1 = (1, 1, 1, 2)$: In this case, it must hold $D_c(X, Y) = c(-, y_1) + D_c(X, Y')$, otherwise $(v_0, v_1) \notin E$. Now, if Y' is empty, then p must have the form $p = (1, 1, 1, 1), (1, 1, 1, 2), (1, 2, 1, 2), \dots, (1, |X| + 1, 1, 2)$, and \emptyset must be a co-optimal mapping between X and Y'. Accordingly, $\emptyset = M_p$ must also be a co-optimal mapping between X and Y, because $c(\emptyset, X, Y) = c(-, y_1) + c(\emptyset, X, Y') = c(-, y_1) + D_c(X, Y') = D_c(X, Y)$.
 - If Y' is not empty, then the path v_1, \ldots, v_T is isomorphic to a path from (1,1,1,1) to (1,|X|+1,1,|Y'|+1) in $\mathcal{G}_{c,X,Y'}$. Accordingly, by virtue of our induction hypothesis, $M_{p'}$ is a co-optimal mapping between X and Y'. Further, we obtain per construction $M_p = \{(i,j+1)|(i,j) \in M_{p'}\}$. Accordingly, it holds: $c(M_p,X,Y) = c(-,y_1)+c(M_{p'},X,Y') = c(-,y_1)+D_c(X,Y') = D_c(X,Y)$, which means that M_p is co-optimal, as claimed.
- $v_1=(1,2,1,2)$: In this case, it must hold $D_c(X,Y)=c(x_1,y_1)+D_c(X',Y')$. Now, if X' is empty, then p must have the form $p=(1,1,1,1), (1,2,1,2), \ldots, (1,2,1,|Y|+1)$, and \emptyset must be a co-optimal mapping between X' and Y'. Accordingly, $\{(1,1)\}=M_p$ must also be a co-optimal mapping between X and Y, because $c(\{(1,1)\},X,Y)=c(x_1,y_1)+c(\emptyset,X',Y')=c(x_1,y_1)+D_c(X',Y')=D_c(X,Y)$.
 - If Y' is empty, then p must have the form $p = (1, 1, 1, 1), (1, 2, 1, 2), \dots, (1, |X| + 1, 1, 2)$, and \emptyset must be a co-optimal mapping between X and Y'. Accordingly, $\{(1,1)\} = M_p$ must also be a co-optimal mapping between X and Y, because $c(\{(1,1)\}, X, Y) = c(x_1, y_1) + c(\emptyset, X', Y') = c(x_1, y_1) + D_c(X', Y') = D_c(X, Y)$.
 - If neither X' nor Y' are empty, then the path v_1, \ldots, v_T is isomorphic to a path from (1, 1, 1, 1) to (1, |X'| + 1, 1, |Y'| + 1) in $\mathcal{G}_{c, X', Y'}$. Accordingly, by virtue of our induction hypothesis, $M_{p'}$ is a co-optimal mapping between X' and Y'. Further, we obtain per construction $M_p = \{(1, 1)\} \cup \{(i+1, j+1) | (i, j) \in M_{p'}\}$. Accordingly, it holds: $c(M_p, X, Y) = c(x_1, y_1) + c(M_{p'}, X', Y') = c(x_1, y_1) + D_c(X', Y') = D_c(X, Y)$, which means that M_p is co-optimal, as claimed.

Other cases can not occur such that our induction is concluded.

Regarding the second claim, let M be a co-optimal mapping between X and Y, i.e. $c(M, X, Y) = D_c(X, Y)$, and distinguish the following cases.

- $1 \in I(M,X,Y)$: In this case it holds $c(M,X,Y) = c(x_1,-) + c(M',X',Y)$ with $M' = \{(i-1,j) | (i,j) \in M\}$. It must hold that M' is a co-optimal mapping between X' and Y. Otherwise, we would obtain $D_c(X,Y) \leq D_c(X',Y) + c(x_1,-) < c(M',X',Y) + c(x_1,-) = c(M,X,Y) = D_c(X,Y)$, which is a contradiction. This also implies that $D_c(X,Y) = c(x_1,-) + D_c(X',Y)$, which in turn implies that $((1,1,1,1),(1,2,1,1)) \in E$.
 - Now, if $X' = \epsilon$, M must be \emptyset , and we can construct the path $p = (1, 1, 1, 1), (1, 2, 1, 1), \dots, (1, 2, 1, |Y| + 1)$, which is a path through $\mathcal{G}_{c,X,Y}$ such that $M_p = \emptyset$.

If X' is not empty, our induction hypothesis implies that there exists a path p' from (1,1,1,1) to (1,|X'|+1,1,|Y|+1) in $\mathcal{G}_{c,X',Y}$ such that $M_{p'}=M'$. Therefore, we can construct an isomorphic path \tilde{p} between (1,2,1,1) and (1,|X|+1,|Y|+1) in $\mathcal{G}_{c,X,Y}$. Accordingly, $p:=(1,1,1,1), \tilde{p}$ must be a path through $\mathcal{G}_{c,X,Y}$, and per construction it must hold that $M_p=M$.

 $1 \in J(M,X,Y)$: In this case it holds $c(M,X,Y) = c(-,y_1) + c(M',X,Y')$ with $M' = \{(i,j-1)|(i,j) \in M\}$. It must hold that M' is a co-optimal mapping between X and Y'. Otherwise, we would obtain $D_c(X,Y) \leq D_c(X,Y') + c(-,y_1) < c(M',X,Y') + c(-,y_1) = c(M,X,Y) = D_c(X,Y)$, which is a contradiction. This also implies that $D_c(X,Y) = c(-,y_1) + D_c(X,Y')$, which in turn implies that $((1,1,1,1),(1,1,1,2)) \in E$.

Now, if $Y' = \epsilon$, M must be \emptyset , and we can construct the path $p = (1, 1, 1, 1), (1, 1, 1, 2), \dots, (1, |X| + 1, 1, 2)$, which is a path through $\mathcal{G}_{c,X,Y}$ such that $M_p = \emptyset$.

If X' is not empty, our induction hypothesis implies that there exists a path p' from (1,1,1,1) to (1,|X|+1,1,|Y'|+1) in $\mathcal{G}_{c,X,Y'}$ such that $M_{p'}=M'$. Therefore, we can construct an isomorphic path \tilde{p} between (1,2,1,1) and (1,|X|+1,|Y|+1) in $\mathcal{G}_{c,X,Y}$. Accordingly, $p:=(1,1,1,1),\tilde{p}$ must be a path through $\mathcal{G}_{c,X,Y}$, and per construction it must hold that $M_p=M$.

 $\exists (1,j), (i,1) \in M$: In this case, i=j, which we can show as follows. Consider the case j>1. In that case, i<1, which is impossible. Similarly, if i>1, it must hold j<1, which is impossible. Therefore i=1 and j=1.

In this case it holds $c(M, X, Y) = c(x_1, y_1) + c(M', X, Y')$ with $M' = \{(i-1, j-1) | (i, j) \in M \setminus \{(1, 1)\}\}$. It must hold that M' is a co-optimal mapping between X' and Y'. Otherwise, we would obtain $D_c(X, Y) \leq D_c(X', Y') + c(x_1, y_1) < c(M', X', Y') + c(x_1, y_1) = c(M, X, Y) = D_c(X, Y)$, which is a contradiction. This also implies that $D_c(X, Y) = c(x_1, y_1) + D_c(X', Y')$, which in turn implies that $((1, 1, 1, 1), (1, 2, 1, 2)) \in E$.

Now, if $X' = \epsilon$, M must be $\{(1,1)\}$, and we can construct the path $p = (1,1,1,1), (1,2,1,2), \ldots$, (1,2,1,|Y|+1), which is a path through $\mathcal{G}_{c,X,Y}$ such that $M_p = \{(1,1)\}$.

If $Y' = \epsilon$, M must be $\{(1,1)\}$, and we can construct the path $p = (1,1,1,1), (1,2,1,2), \ldots, (1,|X|+1,1,2)$, which is a path through $\mathcal{G}_{c,X,Y}$ such that $M_p = \{(1,1)\}$.

If neither X' nor Y' is empty, our induction hypothesis implies that there exists a path p' through $\mathcal{G}_{c,X',Y'}$ such that $M_{p'}=M'$. Therefore, we can construct an isomorphic path \tilde{p} between (1,2,1,2) and (1,|X|+1,|Y|+1) in $\mathcal{G}_{c,X,Y}$. Accordingly, $p:=(1,1,1,1), \tilde{p}$ must be a path through $\mathcal{G}_{c,X,Y}$, and per construction it must hold that $M_p=M$.

As no other cases can occur, this concludes the proof.

Now that we have proven that finding a co-optimal mapping is equivalent to finding a path through the co-optimal edit graph, it is relatively simple to construct an algorithm which identifies one such mapping.

Theorem 16. Given two input trees \bar{x} and \bar{y} as well as a cost function c that is non-negative, self-equal, and conforms to the triangular inequality, Algorithm 6 computes a co-optimal mapping M between \bar{x} and \bar{y} . Further, Algorithm 6 runs in $\mathcal{O}((m+n) \cdot m \cdot n)$ time complexity and $\mathcal{O}(m \cdot n)$ space complexity.

Proof. Algorithm 6 starts at (1,1,1,1) and then travels along the co-optimal edit graph, implicitly constructing it as needed. In particular, lines 5-10 cover the edges defined in Equations 28 and 29, and lines 12-32 cover the edges defined via Equation 30. Lines 34-35 cover the edges defined via Equation 26, and lines 36-38 cover the edges defined via Equation 27. The backwards connections defined via Equations 31, 32, 33 are automatically taken via the keyroot update in line 13.

Further note that Algorithm 6 directly constructs a co-optimal mapping from the path via line 7 which adds a tuple (i, j) to the mapping whenever we use a replacement edge, as suggested by Theorem 15.

Regarding space complexity, we maintain only the matrices D and d from before, which results in $\mathcal{O}(m \cdot n)$ space complexity. Regarding runtime, we note that the main while loop of Algorithm 6

runs at most $\mathcal{O}(m+n)$ times because in each loop iteration, i or j (or both) are increased, or k and l are updated, such that i or j (or both) are increased in the next iteration. Within each iteration, the worst case is that we need to update a section of \mathbf{D} . Each of these updates takes at worst $\mathcal{O}(m \cdot n)$ operations, such that we obtain $\mathcal{O}((m+n) \cdot m \cdot n)$ overall.

As an example, consider the co-optimal edit graph in Figure 9, corresponding to the subforest edit distances in Figure 7. We begin by considering the entire trees $\bar{x}=a(b(c,d),e)$ and $\bar{y}=f(g)$ by setting i=j=k=l=1. Our initial mapping is empty, i.e. $M=\emptyset$. In this situation we observe that $5=D_{1,1}=D_{2,2}+c(x_1,y_1)=4+1$, that is: replacing a with f is part of a co-optimal mapping. Therefore, we add (1,1) to the mapping and increment both i and j. Now, we find that $4=D_{2,2}=D_{5,3}+d_{2,2}=1+3$, that is: replacing the subtree b(c,d) with the subtree g is part of a co-optimal mapping. Therefore, we update D with the entries $D_{i,j}=D_c(\bar{x}[i,rl_{\bar{x}}(2)],\bar{y}[j,rl_{\bar{y}}(2)])+D_c([rl_{\bar{x}}(2)+1,5],[rl_{\bar{y}}(2)+1,2]=D_c(\bar{x}[i,4],\bar{y}[j,2])+D_c(e,\epsilon)$ for all $i\in\{2,3,4,5\}$ and $j\in\{2,3\}$. In this example, this only changes the entry $D_{5,2}$ which would be set to $D_c(\bar{x}[i,4],\bar{y}[j,2)+D_c(e,\epsilon)=1+1=2$. Notice that we also update the current keyroots to k=2 and l=2.

In the next iteration, we find that $4 = \mathbf{D}_{2,2} = \mathbf{D}_{3,3} + c(x_2, y_2) = 3 + 1$, that is: replacing b with g is part of a co-optimal mapping. Therefore, we add (2,2) to the mapping increment both i and j. Now, 3 = j > n = 2. Therefore, the algorithm stops and returns the mapping $M = \{(1,1), (2,2)\}$ which is indeed a co-optimal mapping between \bar{x} and \bar{y} .

Note that Algorithm 6 always prefers replacements if multiple edits are co-optimal. As such, Algorithm 6 will prefer to map the nodes close to the root of both trees to each other, and delete/insert nodes closer to the leaves. The other possible co-optimal mappings for this example are $\{(1,1),(3,2)\}$, $\{(1,1),(4,2)\}$, $\{(1,1),(5,2)\}$, $\{(2,1),(3,2)\}$, and $\{(2,1),(4,2)\}$ (also refer to Figure 8), all of which correspond to exactly one path in Figure 9.

As we have already seen, it is infeasible to list all co-optimal mappings in general (see Theorem 13). Interestingly, though, we can still count the number of such mappings efficiently. We will first consider the problem of counting the number of paths in a general DAG, and then return to the co-optimal edit graph specifically.

Theorem 17. Let $\mathcal{G} = (V, E)$ be a DAG with ordering relation < and let v_1, \ldots, v_n be the nodes in V as ordered according to <. Then, Algorithm 7 returns a $n \times 1$ vector $\vec{\alpha}$, such that α_i is exactly the number of paths leading from v_1 to v_i . Further, Algorithm 7 runs in $\mathcal{O}(n)$ time and space complexity.

Proof. To prove this result, we first show two lemmata:

- 1. Algorithm 7 visits all reachable nodes from v_1 in ascending order, and no other nodes.
- 2. When Algorithm 7 visits node v_i , α_i contains exactly the number of paths from v_1 to v_i .

We call a note *visited*, if it is pulled from Q. We proof both lemmata by induction over i.

1. Our base case is v_1 , which is indeed visited first.

Now, assume that the claim holds for all reachable nodes $\leq v$. Consider the smallest node u > v which is reachable from 1. Then, there is a path u_0, \ldots, u_T with $u_0 = v_1$ and $u_T = u$. Because \mathcal{G} is a DAG, $u_{T-1} < u$. Further, because u_0, \ldots, u_{T-1} is a path from v_1 to u_{T-1}, u_{T-1} is reachable from v_1 . Because u is per definition the smallest node larger than v which is reachable from v_1 , it must hold $u_{T-1} \leq v$. Therefore, per induction, u_{T-1} has been visited before. This implies that $u \in Q$. Because we select the minimum from Q in each iteration, and because all elements smaller than u have been visited before (and are not visited again due to the DAG property), u will be visited next. Therefore, still all reachable nodes from v_1 are visited in ascending order, and all nodes that are visited are reachable nodes.

2. Again, our base case is v_1 , which is visited first. As it is visited, $\alpha_1 = 1$. Indeed, there is only one path from v_1 to v_1 , which is the trivial path.

Now, assume that the claim holds for all reachable nodes $\leq v$. Then, consider the smallest node $v_i > v$ which is reachable from v_1 . Further, let v_{i_1}, \ldots, v_{i_m} be all nodes which are reachable from v_1 , such that $(v_{i_j}, v_i) \in E$. Because \mathcal{G} is a DAG, $v_{i_j} < v_i$. Further, because v_{i_j} is reachable from 1 and v_i is per definition the smallest node larger than v which is reachable from v_1 , it

Algorithm 6 A backtracing algorithm for the TED, which infers a co-optimal mapping between the input trees \bar{x} and \bar{y} . Refer to our project web site for a reference implementation.

```
1: function BACKTRACE (Two trees \bar{x} and \bar{y}, the matrices d and D after executing Algorithm 5, and
     a cost function c)
          M \leftarrow \emptyset.
 2:
 3:
          i \leftarrow 1, j \leftarrow 1, k \leftarrow 1, l \leftarrow 1.
 4:
          while i \leq m \land j \leq n do
               if (rl_{\bar{x}}(i) = rl_{\bar{x}}(k) \wedge rl_{\bar{y}}(j) = rl_{\bar{y}}(l)) \vee (c(x_i, y_j) = c(x_i, -) + c(-, y_j)) then
 5:
                    if D_{i,j} = D_{i+1,j+1} + c(x_i, y_j) then
 6:
                         M \leftarrow M \cup \{(i,j)\}.
                                                                                                                > replacement is optimal
 7:
                         i \leftarrow i + 1, j \leftarrow j + 1.
 8:
                         continue.
 9:
                    end if
10:
               else
11:
12:
                    if D_{i,j} = D_{rl_{\bar{x}}(i)+1, rl_{\bar{y}}(j)+1} + d_{i,j} then
                         k \leftarrow i, l \leftarrow j.
                                                                                   \triangleright subtree replacement is optimal; update D.
13:
                         for i' \leftarrow rl_{\bar{x}}(k), \ldots, k do
14:
                               D_{i',rl_{\bar{y}}(l)+1} \leftarrow D_{i'+1,rl_{\bar{y}}(l)+1} + c(x_{i'},-).
                                                                                                                                ⊳ Equation 21
15:
                         end for
16:
                         for j' \leftarrow rl_{\bar{y}}(l), \ldots, l do
17:
                              D_{rl_{\bar{x}}(k)+1,j'} \leftarrow D_{rl_{\bar{x}}(k)+1,j'+1} + c(-,y_{j'}).
                                                                                                                                ⊳ Equation 22
18:
19:
                         for i' \leftarrow rl_{\bar{x}}(k), \ldots, k do
20:
                               for j' \leftarrow rl_{\bar{y}}(l), \ldots, l do
21:
                                   if rl_{\bar{x}}(i') = rl_{\bar{x}}(k) \wedge rl_{\bar{y}}(j') = rl_{\bar{y}}(l) then
22:
23:
                                         D_{i',j'} \leftarrow d_{i',j'} + D_{rl_{\bar{x}}(k)+1,rl_{\bar{y}}(l)+1}.
24:
                                   else
                                         D_{i',j'} \leftarrow \min\{D_{i'+1,j'} + c(x_{i'}, -),
25:
                                            D_{i',j'+1} + c(-,y_{j'}),
26:
                                            D_{rl_{\bar{x}}(i')+1,rl_{\bar{y}}(j')+1}+d_{i',j'}\}.
                                                                                                                                ⊳ Equation 15
27:
                                    end if
28:
                               end for
29:
                         end for
30:
                         continue.
31:
                    end if
32:
33:
               end if
               if D_{i,j} = D_{i+1,j} + c(x_i, -) then
34:
                    i \leftarrow i+1.
                                                                                                                      \triangleright deletion is optimal
35:
               else if D_{i,j} = D_{i,j+1} + c(-, y_j) then
36:
37:
                    j \leftarrow j + 1.
                                                                                                                     end if
38:
          end while
39:
          return M.
40:
41: end function
```

must hold $v_{i_j} \leq v$. Therefore, per induction, α_{i_j} is equal to the number of paths from v_1 to v_{i_j} . For any such path p, the concatenation $p \oplus v_i$ is a path from v_1 to v_i . Conversely, we can decompose any path p' from v_1 to v_i as $p' = p \oplus v_i$ where p is a path from v_1 to some node v_{i_j} . Accordingly, the number of paths from v_1 to v_i is exactly $\sum_{j=1}^m \alpha_{i_j}$.

Finally, because of the first lemma, we know that all v_{i_j} have been visited already (without duplicates), and that on each of these visits, α_{i_j} has been added to α_i . Therefore, we obtain $\alpha_i = \sum_{j=1}^m \alpha_{i_j}$.

Because Lemma 1 implies that we do not visit any node smaller than v_i after v_i has been visited, the value α_i does not change after v_i is visited. Therefore, α_i still contains the number of paths from v_1 to v_i at the end of the algorithm.

Regarding runtime, it follows from the first lemma that, per iteration, exactly one reachable node is processed and will not be visited again. In the worst case, all nodes in the graph are reachable, which yields $\mathcal{O}(n)$ iterations. In each iteration we need to retrieve the minimum of Q and insert all v into Q, for which $(u,v) \in E$. Both is possible in constant time if a suitable data structure for Q is used. If one uses a tree structure for Q, the runtime rises to $\mathcal{O}(n \cdot \log(n))$. The space complexity is $\mathcal{O}(n)$ because $\vec{\alpha}$ has n entries and Q can not exceed n entries.

```
Algorithm 7 An algorithm to count the number of paths between v_1 and v_i in a DAG \mathcal{G} = (\{v_1, \ldots, v_n\}, E) with ordering relation <.
```

```
function COUNT-PATHS-FORWARD(A DAG \mathcal{G} = (\{v_1, \dots, v_n\}, E) with ordering relation <.) \vec{\alpha} \leftarrow \vec{0}^n. \alpha_1 \leftarrow 1. Q \leftarrow \{v_1\}. while Q \neq \emptyset do v_i \leftarrow \min_{<} Q. Q \leftarrow Q \setminus \{v_i\}. for (v_i, v_j) \in E do \alpha_j \leftarrow \alpha_j + \alpha_i. Q \leftarrow Q \cup \{v_j\}. end for end while return \vec{\alpha}. end function
```

As an example, consider the DAG in Figure 9 with the ordering indices shown in blue. Assuming a sorted set for Q, Algorithm 7 would initialize $\alpha_1 \leftarrow 1$ and $Q \leftarrow \{(1,1,1,1)\}$ and would then behave as follows.

```
1. v_i = v_1 = (1, 1, 1, 1), \alpha_2 \leftarrow 1, \alpha_3 \leftarrow 1, Q \leftarrow \{(1, 2, 1, 1), (1, 2, 1, 2)\}.

2. v_i = v_2 = (1, 2, 1, 1), \alpha_4 \leftarrow 1, Q \leftarrow \{(1, 2, 1, 2), (2, 3, 1, 2)\}.

3. v_i = v_3 = (1, 2, 1, 2), \alpha_5 \leftarrow 1, \alpha_6 \leftarrow 1, Q \leftarrow \{(2, 3, 1, 2), (1, 3, 1, 2), (2, 3, 1, 3)\}.

4. v_i = v_4 = (2, 3, 1, 2), \alpha_7 \leftarrow 1, \alpha_9 \leftarrow 1, Q \leftarrow \{(1, 3, 1, 2), (2, 3, 1, 3), (2, 4, 1, 2), (3, 4, 1, 3)\}.

5. v_i = v_5 = (1, 3, 1, 2), \alpha_8 \leftarrow 1, \alpha_9 \leftarrow 1 + 1, Q \leftarrow \{(2, 3, 1, 3), (2, 4, 1, 2), (1, 4, 1, 2), (3, 4, 1, 3)\}.

6. v_i = v_6 = (2, 3, 1, 3), \alpha_{10} \leftarrow 1, Q \leftarrow \{(2, 4, 1, 2), (1, 4, 1, 2), (3, 4, 1, 3), (2, 4, 1, 3)\}.

7. v_i = v_7 = (2, 4, 1, 2), \alpha_{12} \leftarrow 1, Q \leftarrow \{(1, 4, 1, 2), (3, 4, 1, 3), (2, 4, 1, 3), (2, 5, 1, 3)\}.

8. v_i = v_8 = (1, 4, 1, 2), \alpha_{11} \leftarrow 1, \alpha_{13} \leftarrow 1, Q \leftarrow \{(3, 4, 1, 3), (2, 4, 1, 3), (1, 5, 1, 2), (2, 5, 1, 3), (1, 5, 1, 3)\}.

9. v_i = v_9 = (3, 4, 1, 3), \alpha_{10} \leftarrow 1 + 2, Q \leftarrow \{(2, 4, 1, 3), (1, 5, 1, 2), (2, 5, 1, 3), (1, 5, 1, 3)\}.

10. v_i = v_{10} = (2, 4, 1, 3), \alpha_{12} \leftarrow 1 + 3, Q \leftarrow \{(1, 5, 1, 2), (2, 5, 1, 3), (1, 5, 1, 3)\}.
```

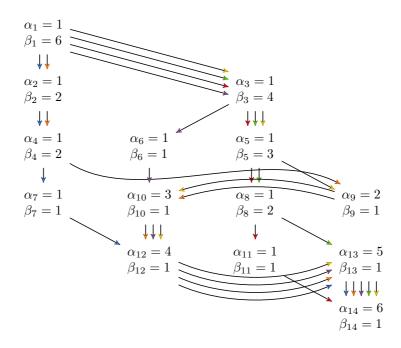


Figure 10: An illustration of all possible paths from (1,1,1,1) to (1,6,1,3) in the DAG from Figure 9. The paths are differentiated by color. The number of (unique) paths from (1,1,1,1) to each node v_i is denoted as α_i , and the number of (unique) paths from each node v_i to (1,6,1,3) is denoted as β_i .

11.
$$v_i = v_{11} = (1, 5, 1, 2), \alpha_{14} \leftarrow 1, Q \leftarrow \{(2, 5, 1, 3), (1, 5, 1, 3), (1, 6, 1, 3)\}.$$

12.
$$v_i = v_{12} = (2, 5, 1, 3), \alpha_{13} \leftarrow 1 + 4, Q \leftarrow \{(1, 5, 1, 3), (1, 6, 1, 3)\}.$$

13.
$$v_i = v_{13} = (1, 5, 1, 3), \alpha_{14} \leftarrow 1 + 5, Q \leftarrow \{(1, 6, 1, 3)\}.$$

14.
$$v_i = v_{14} = (1, 6, 1, 3), Q \leftarrow \emptyset.$$

The resulting α -values for all nodes are indicated in Figure 10.

Interestingly, we can also invert this computation to compute the number of paths which lead from any node v to the last node in a DAG.

Theorem 18. Let $\mathcal{G} = (V, E)$ be a DAG with ordering relation < and let v_1, \ldots, v_n be the nodes in V as ordered according to <. Then, Algorithm 8 returns a $n \times 1$ vector $\vec{\beta}$, such that β_i is exactly the number of paths leading from v_i to v_n . Further, Algorithm 8 runs in $\mathcal{O}(n)$ time and space complexity.

Proof. Note that the structure of this proof is exactly symmetric to Theorem 17.

To prove this result, we first show two lemmata:

- 1. Algorithm 8 visits all nodes from which v_n is reachable in descending order, and no other nodes.
- 2. When Algorithm 8 visits node v_i , β_i contains exactly the number of paths from v_i to v_n .

We call a note visited, if it is pulled from Q. We proof both lemmata by induction over i in descending order.

1. Our base case is v_n , which is indeed visited first.

Now, assume that the claim holds for nodes $\geq v$ such that v_n is reachable from v. Consider now the largest v_i , such that $v > v_i$ and v_n is reachable from v_i . Then, there is a path u_0, \ldots, u_T with $u_0 = v_i$ and $u_T = v_n$. Because \mathcal{G} is a DAG, $u_1 > v_j$. Further, because u_1, \ldots, u_T is a path from u_1 to v_n , v_n is reachable from u_1 . Because v_i is per definition the largest node smaller than v from which v_n is reachable, $u_1 \geq v$. Therefore, per induction, u_1 has been visited before. This

implies that $v_i \in Q$. Because we select the maximum from Q in each iteration, and because all elements larger than v_i have been visited before (and are not visited again due to the DAG property), v_i will be visited next. Therefore, still all nodes from which v_n is reachable are visited in descending order, and all nodes which are visited are nodes from which v_n is reachable.

2. Again, our base case is v_n , which is visited first. As it is visited, we have $\beta_n = 1$. And indeed there is only one path from v_n to v_n , namely the trivial path.

Now, assume that the claim holds for all nodes $\geq v$ from which v_n is reachable. Then, consider the largest node $v_i < v$ from which v_n is reachable. Further, let v_{i_1}, \ldots, v_{i_m} be all nodes from which v_n is reachable, such that $(v_i, v_{i_j}) \in E$. Because $\mathcal G$ is a DAG, $v_{i_j} > v_i$ for all j. Further, because v_n is reachable from v_{i_j} and v_i is the largest node smaller than v from which v_n is reachable, it must hold $v_{i_j} \geq v$. Therefore, per induction, β_{i_j} is the number of paths from v_{i_j} to v_n . For any such path p, the concatenation $v_i \oplus p$ is a path from v_i to v_n . Conversely, we can decompose any path p' from v_i to v_n as $p' = v_i \oplus p$ where p is a path from v_{i_j} to v_n for some j. Accordingly, the number of paths from v_i to v_n is exactly $\sum_{j=1}^m \beta_{i_j}$.

Finally, because of the first lemma, we know that all v_{i_j} have been visited already (without duplicates), and that on each of these visits, β_{i_j} has been added to β_i . Therefore, we obtain $\beta_i = \sum_{j=1}^m \beta_{i_j}$.

Because Lemma 1 implies that we do not visit any node larger than v_i after v_i has been visited, the value β_i does not change after i is visited. Therefore, β_i still contains the number of paths from v_i to v_T at the end of the algorithm.

Regarding runtime, it follows from the first lemma that, per iteration, exactly one reachable node is processed and will not be visited again. In the worst case, all nodes in the graph are reachable, which yields $\mathcal{O}(n)$ iterations. In each iteration we need to retrieve the maximum of Q and insert all u into Q, for which $(u, v) \in E$. Both is possible in constant time if a suitable data structure for Q is used. If one uses a tree structure for Q, the runtime rises to $\mathcal{O}(n \cdot \log(n))$. The space complexity is $\mathcal{O}(n)$ because $\vec{\beta}$ has n entries and Q can not exceed n entries.

Algorithm 8 An algorithm to count the number of paths between each node v_i and node v_n in a DAG $\mathcal{G} = (\{v_1, \ldots, v_n\}, E)$, where v_1, \ldots, v_n is the ordered node list according to the DAGs ordering relation <.

```
function COUNT-PATHS-BACKWARD(A DAG \mathcal{G} = (\{v_1, \dots, v_n\}, E), an ordering relation <.)
\vec{\beta} \leftarrow n \times 1 \text{ vector of zeros.}
\beta_n \leftarrow 1.
Q \leftarrow \{v_n\}.
\text{while } Q \neq \emptyset \text{ do}
v_j \leftarrow \max_{<} Q.
Q \leftarrow Q \setminus \{v_j\}.
\text{for } (v_i, v_j) \in E \text{ do}
\beta_i \leftarrow \beta_i + \beta_j.
Q \leftarrow Q \cup \{v_i\}.
\text{end for}
\text{end while}
\text{return } \vec{\beta}.
end function
```

As an example, consider the DAG on the right in Figure 9. The resulting β -values for all nodes are indicated in Figure 10.

Beyond the utility of counting the number of paths in linear time, the combination of both algorithms also permits us to compute how often a certain edge of the graph occurs in paths from v_1 to v_n .

Theorem 19. Let $\mathcal{G} = (V = \{v_1, \dots, v_n\}, E)$ be a DAG with ordering relation < where v_1, \dots, v_n are ordered according to <. Further, let $\vec{\alpha}$ be the result of Algorithm 7 for \mathcal{G} , and let $\vec{\beta}$ be the result of

Algorithm 8 for \mathcal{G} . Then, for any edge $(v_i, v_j) \in E$ it holds: $\alpha_i \cdot \beta_j$ is precisely the number of paths from v_1 to v_n which contain v_i, v_j , that is, paths $p = u_1, \ldots, u_T$ such that an $t \in \{1, \ldots, T-1\}$ exists for which $u_t = v_i$ and $u_{t+1} = v_j$.

Proof. Let $(v_i, v_j) \in E$ and let m be the number of paths from v_1 to v_n which traverse (u, v). Further, let $u_1, \ldots, u_{T'}$ be a path from v_1 to v_i and let $u_{T'+1}, \ldots, u_T$ be a path from v_j to v_n . Then, because $(v_i, v_j) \in E, u_1, \ldots, u_T$ is a path from v_1 to v_n in $\mathcal G$ which traverses (v_i, v_j) . We know by virtue of Theorem 17 that the number of paths from v_1 to v_i is α_i , and we know by virtue of Theorem 18 that the number of paths from v_j to v_n is β_j . Now, as we noted before, any combination of a path counted in α_i and a path counted in β_j is a path from v_1 to v_n , and any of these combinations is unique. Therefore, we obtain $m \geq \alpha_u \cdot \beta_v$.

Further, we note that we can decompose any path from v_1 to v_n as illustrated above, such that $m \leq \alpha_u \cdot \beta_v$.

For the example DAG from Figure 9 we show all possible paths from (1,1,1,1) to (1,6,1,3) in Figure 10. For every edge in this DAG you can verify that, indeed, the number of traversing paths is equivalent to the α value of the source node times the β value of the target node.

Note that the number of paths which traverses a certain edge reveals crucial information about the co-optimal mappings. In particular, if we consider an edge of the form ((k, i, l, j), (k', i + 1, l', j + 1)), the number of paths from (1, 1, 1, 1) to (1, |X| + 1, 1, |Y| + 1) which traverse this edge is an estimate of the number of co-optimal mappings which contain the tuple (i, j). Unfortunately, this estimate is not necessarily exact, because there may be multiple paths through the co-optimal edit graph which correspond to the same co-optimal mapping.

In particular, excessive paths occure whenever $c(x_i, -) + c(-, y_j) = c(x_i, y_j)$. In these cases, deletion, replacement, and insertion are all co-optimal, and thus there exist three paths from (k, i, l, j) to (k, i + 1, l, j + 1), one which uses a deletion first and then an insertion, one which uses only a replacement, and one which uses an insertion first and then a deletion. The first and last of these paths correspond to the same co-optimal mapping, leading to overcounting. We can avoid this overcounting by covering this special case explicitly. This results in a new forward-counting-Algorithm 9, a new backward-counting-Algorithm 10, and a forward-backward Algorithm 11 which characterize the number of co-optimal mappings rather than the number of co-optimal paths.

Theorem 20. Let \bar{x} and \bar{y} be trees over some alphabet \mathcal{X} and let c be a cost function that is nonnegative, self-equal, and conforms to the triangular inequality. Further, let $\mathcal{G}_{\bar{x},\bar{y},c} = (V,E)$ be the co-optimal edit graph corresponding to \bar{x} , \bar{y} , and c. Then, the first output argument of Algorithm 11 is a $|\bar{x}| \times |\bar{y}|$ matrix Γ such that $\Gamma_{i,j}$ is exactly the number of co-optimal mappings which contain (i,j). Further, the second output argument of Algorithm 11 is the number of co-optimal mappings. Finally, Algorithm 11 has $\mathcal{O}(|\bar{x}|^6 \cdot |\bar{y}|^6)$ time and $\mathcal{O}(|\bar{x}|^2 \cdot |\bar{y}|^2)$ space complexity.

Proof. For the technical details of this proof, refer to my dissertation (Paaßen 2018). Here, I provide a sketch of the proof.

First, we observe that Algorithm 9 is analogous to Algorithm 7, and that Algorithm 10 is analogous to Algorithm 8. The latter analogy holds because we just postpone adding the contributions to β_i to the visit of β_i itself, but all contributions are still collected. We further speed up the process by considering only cells of the dynamic programming matrix which are actually reachable from (1,1,1,1). Another non-obvious part of the analogy is that we go into recursion to compute the number of cooptimal paths for a subtree replacement. In this regard, we note that we can extend each path from (1,1,1,1) to (k,i,l,j) to a path to $(k,rl_{\bar{x}}(i)+1,l,rl_{\bar{y}}(j)+1)$ by using one of the possible paths in the co-optimal edit graph corresponding to \bar{x}_i and \bar{y}_j . However, this would over-count the paths which delete the node x_i or insert the node y_j , which we prevent by setting $\mathbf{D}'_{1,2} = \mathbf{D}'_{2,1} = \infty$. The same argument holds for the backwards case: We can extend any path from $(k, rl_{\bar{x}}(i)+1, l, rl_{\bar{y}}(j)+1)$ to $(1, |\bar{x}|+1, 1, |\bar{y}|+1)$ to a path from (i,j) to $(|\bar{x}|+1, |\bar{y}|+1)$ by using one of the possible paths in the co-optimal edit graph corresponding to \bar{x}_i and \bar{y}_j .

Finally, algorithm 11 computes the products of α and β -values according to Theorem 19. The only special case is, once again, the case of subtree replacements. In that case, we can again argue that, for any combination of a path which leads from (1,1,1,1) to (k,i,l,j), and a path which leads from $(k,rl_{\bar{x}}(i)+1,l,rl_{\bar{y}}(j)+1)$ to $(1,|\bar{x}|+1,1,|\bar{y}|+1)$, we can construct a path from (1,1,1,1) to

Table 1: The forward matrix A, the backward matrix B, and the matrix Γ for the trees $\bar{x} = a(b(c,d),e)$ and $\bar{y} = f(g)$ from Figure 7, as returned by Algorithms 9, 10, and 11 respectively. The color coding follows Figure 7.

4	á	1 2	3	\mathbf{p}	٠	i	1	2	3
$egin{aligned} oldsymbol{A}_{i,j} \ i \end{aligned}$	$x_i \setminus y_j$	f g	- -	D	$i^{i,j}$	$x_i \setminus y_j$	f	g	- -
1	a	1			1	a	16 K		
2	b	41/41/		6	2	b	125	4	
3	С	1/1	\		3	С	,/	/34	\
4	d	ì	1	4	4	d		2	1
5	е	1	¥5¢		5	е		1	1
6	_		164	(6	_		,	\backslash_1
		$egin{array}{c} oldsymbol{\Gamma}_{i,j} \ i \end{array}$	j	1	2				
		i	$x_i \setminus y_j$	f	g				
		1	a	4	0				
		2	b	2	1				
		3	C	Ω	1 +	1			

 $(1, |\bar{x}|+1, 1, |\bar{y}|+1)$ by inserting a "middle piece" which corresponds to a path in the co-optimal edit graph for \bar{x}_i and \bar{y}_j . Therefore, for $\gamma = A_{i,j} \cdot B_{rl_{\bar{x}}(i)+1, rl_{\bar{y}}(j)+1}$, $\gamma \cdot \Gamma'_{i',j'}$ is an additional contribution to the count of (i'+i-1, j'+j-1).

d

0

1 + 1

4

Now, consider the efficiency claims. First, we analyze Algorithm 7. In the worst case, lines 27-31 need to be executed in each possible iteration. In that case, \mathbf{D}' and \mathbf{d}' need to be computed via Algorithm 5, which requires $\mathcal{O}(|\bar{x}|^2 \cdot |\bar{y}|^2)$ steps and $\mathcal{O}(|\bar{x}| \cdot |\bar{y}|)$ space. Including the recursive calls, this can occur $\mathcal{O}(|\bar{x}|^2 \cdot |\bar{y}|^2)$ times at worst such that Algorithm 7 has an overall runtime complexity of $\mathcal{O}(|\bar{x}|^4 \cdot |\bar{y}|^4)$.

Regarding space complexity, each level of recursion needs to maintain a constant number of matrices of size $\mathcal{O}(|\bar{x}| \cdot |\bar{y}|)$. A worst, there can be $\mathcal{O}(|\bar{x}| \cdot |\bar{y}|)$ levels of recursion active at the same time, implying a space complexity of $\mathcal{O}(|\bar{x}|^2 \cdot |\bar{y}|^2)$.

Now, note that Algorithm 8, by construction, iterates over the same elements as Algorithm 7 and has the same structure, such that the complexity results carry over.

Finally, regarding Algorithm 9 itself, we find that, in the worst case, lines 15-23 get executed in every possible iteration. These lines include a recursive call to Algorithm 9, and in each such recursive call, Algorithm 7 and Algorithm 8 get executed. With the same argument as before, we perform at most $\mathcal{O}(|\bar{x}|^2 \cdot |\bar{y}|^2)$ of such recursive calls, yielding an overall runtime complexity of $\mathcal{O}(|\bar{x}|^6 \cdot |\bar{y}|^6)$ in the worst case.

Regarding space complexity, each level of recursion needs to maintain a constant number of matrices of size $\mathcal{O}(|\bar{x}| \cdot |\bar{y}|)$. A worst, there can be $\mathcal{O}(|\bar{x}| \cdot |\bar{y}|)$ levels of recursion active at the same time, implying a space complexity of $\mathcal{O}(|\bar{x}|^2 \cdot |\bar{y}|^2)$.

Note that the version of the algorithm presented here is dedicated to minimize space complexity. By additionally tabulating Γ for all subtrees, space complexity rises to $\mathcal{O}(|\bar{x}|^4 \cdot |\bar{y}|^4)$ in the worst case, but runtime complexity is reduced to $\mathcal{O}(|\bar{x}|^3 \cdot |\bar{y}|^3)^1$. Another point to note is that the worst case for this algorithm is quite unlikely. First, both input trees would have to be left- or right-heavy. Second, in every step of the computation, multiple options have to be co-optimal, which only occurs in degenerate cases where, for example, the deletion or insertion cost for all symbols is zero.

For the example TED calculation in Figure 7, the results for A, B, and Γ according to Algorithm 11 are shown in Table 4. By comparing with the co-optimal mappings in Figure 8, you can verify that this matrix does indeed sum up all co-optimal mappings.

 $^{^{1}\}mathrm{This}$ is the version we implemented in our reference implementation

Algorithm 9 A variation of the forward path-counting Algorithm 7 for the TED.

```
1: function FORWARD (Two trees \bar{x} and \bar{y}, the matrices d and D after executing Algorithm 5, and a
     cost function c
           Initialize A as a (|\bar{x}|+1) \times (|\bar{y}|+1) matrix of zeros.
 2:
 3:
           A_{1,1} \leftarrow 1, Q \leftarrow \{(1,1)\}
 4:
           while Q \neq \emptyset do
 5:
                (i,j) \leftarrow \min Q.
                                                                                                                        ▶ Lexicographic ordering
 6:
 7:
                Q \leftarrow Q \setminus \{(i,j)\}.
                C \leftarrow C \cup \{(i,j)\}.
 8:
                if i \le |\bar{x}| \wedge D_{i,j} = c(x_i, -) + D_{i+1,j} then
 9:
                      A_{i+1,j} \leftarrow A_{i+1,j} + A_{i,j}.
10:
                      Q \leftarrow Q \cup \{(i+1,j)\}.
11:
12:
                end if
                if j \le |\bar{y}| \land D_{i,j} = c(-,y_j) + D_{i,j+1} then
13:
14:
                      A_{i,j+1} \leftarrow A_{i,j+1} + A_{i,j}.
                      Q \leftarrow Q \cup \{(i, j+1)\}.
15:
16:
                if i = |\bar{x}| + 1 \lor j = |\bar{y}| + 1 \lor c(x_i, y_j) = c(x_i, -) + c(-, y_j) then
17:
                      continue
18:
19:
                end if
                if rl_{\bar{x}}(i) = rl_{\bar{x}}(1) \wedge rl_{\bar{y}}(j) = rl_{\bar{y}}(1) then
20:
                      if D_{i,j} = D_{i+1,j+1} + c(x_i, y_j) then
21:
                           oldsymbol{A}_{i+1,j+1} \leftarrow oldsymbol{A}_{i+1,j+1} + oldsymbol{A}_{i,j}
22:
                           Q \leftarrow Q \cup \{(i+1,j+1)\}.
23:
                      end if
24:
                else
25:
                      if oldsymbol{D}_{i,j} = oldsymbol{D}_{rl_{ar{x}}(i)+1,rl_{ar{y}}(j)+1} + oldsymbol{d}_{i,j} then
26:
                           Compute D' and d' via Algorithm 5 for the subtrees \bar{x}^i and \bar{y}^j.
27:
28:
                           D'_{1,2} \leftarrow \infty. \ D'_{2,1} \leftarrow \infty.
                           (Q', \mathbf{A}') \leftarrow \text{FORWARD}(\bar{x}^i, \bar{y}^j, \mathbf{d}', \mathbf{D}', c).
29:
                           A_{rl_{\bar{x}}(i)+1,rl_{\bar{y}}(j)+1} \leftarrow A_{rl_{\bar{x}}(i)+1,rl_{\bar{y}}(j)+1} + A'_{|\bar{x}^i|+1,|\bar{y}^j|+1} \cdot A_{i,j}.
30:
                           Q \leftarrow Q \cup \{(rl_{\bar{x}}(i) + 1, rl_{\bar{y}}(j) + 1)\}.
31:
                      end if
32:
                end if
33:
           end while
34:
           return (C, \mathbf{A}).
35:
36: end function
```

Algorithm 10 A variation of the backward path-counting Algorithm 8 for the TED.

```
1: function BACKWARD (Two trees \bar{x} and \bar{y}, the matrices d and D after executing Algorithm 5, a
      cost function c, and a set of tuples C as returned by Algorithm 9.)
            Initialize {\pmb B} as a (|\bar x|+1)\times(|\bar y|+1) matrix of zeros.
 2:
 3:
             B_{|\bar{x}|+1,|\bar{y}|+1} \leftarrow 1.
             while C \neq \emptyset do
 4:
                   (i,j) \leftarrow \max C.
                                                                                                                                          ▶ Lexicographic ordering
 5:
                  C \leftarrow C \setminus \{(i,j)\}.
 6:
                  if i \leq |\bar{x}| \wedge D_{i,j} = c(x_i, -) + D_{i+1,j} then
 7:
                         \boldsymbol{B}_{i,j} \leftarrow \boldsymbol{B}_{i,j} + \boldsymbol{B}_{i+1,j}
 8:
 9:
                   if j \le |\bar{y}| \land D_{i,j} = c(-,y_j) + D_{i,j+1} then
10:
                         \boldsymbol{B}_{i,j} \leftarrow \boldsymbol{B}_{i,j} + \boldsymbol{B}_{i,j+1}
11:
12:
                   if i = |\bar{x}| + 1 \lor j = |\bar{y}| + 1 \lor c(x_i, y_j) = c(x_i, -) + c(-, y_j) then
13:
14:
                         continue
15:
                   end if
                   if rl_{\bar{x}}(i) = rl_{\bar{x}}(1) \wedge rl_{\bar{y}}(j) = rl_{\bar{y}}(1) then
16:
                         if D_{i,j} = D_{i+1,j+1} + c(x_i, y_j) then
17:
                                \boldsymbol{B}_{i,j} \leftarrow \boldsymbol{B}_{i,j} + \boldsymbol{B}_{i+1,j+1}
18:
19:
                         end if
                   else
20:
                         \begin{array}{l} \text{if } \boldsymbol{D}_{i,j} = \boldsymbol{D}_{rl_{\bar{x}}(i)+1,rl_{\bar{y}}(j)+1} + \boldsymbol{d}_{i,j} \text{ then} \\ \text{Compute } \boldsymbol{D}' \text{ and } \boldsymbol{d}' \text{ via Algorithm 5 for the subtrees } \bar{x}^i \text{ and } \bar{y}^j. \end{array}
21:
22:
                               D'_{1,2} \leftarrow \infty. \ D'_{2,1} \leftarrow \infty.
23:
                               (Q', \mathbf{A}') \leftarrow \text{FORWARD}(\bar{x}^i, \bar{y}^j, \mathbf{d}', \mathbf{D}', c).
24:
                               B_{i,j} \leftarrow B_{i,j} + B_{rl_{\bar{x}}(i)+1,rl_{\bar{y}}(j)+1} \cdot A'_{|\bar{x}^i|+1,|\bar{y}^j|+1}.
25:
                         end if
26:
                   end if
27:
            end while
28:
            return B.
29:
30: end function
```

Algorithm 11 A forward-backward algorithm to compute the number of times the tuple (i, j) occurs in co-optimal mappings for paths in the co-optimal edit graphs between two input trees \bar{x} and \bar{y} . The second output is the overall number of co-optimal mappings. Refer to our project web site for a reference implementation.

```
1: function COOPTIMALS(Two trees \bar{x} and \bar{y}, the matrices d and D after executing algorithm 5, and
     a cost function c)
           (C, \mathbf{A}) \leftarrow \text{FORWARD}(\bar{x}, \bar{y}, \mathbf{d}, \mathbf{D}, c).
                                                                                                                               \triangleright Refer to Algorithm 9.
 2:
           \boldsymbol{B} \leftarrow \text{BACKWARD}(\bar{x}, \bar{y}, \boldsymbol{d}, \boldsymbol{D}, c, C).
                                                                                                                             ▶ Refer to Algorithm 10.
 3:
           Initialize \Gamma as a |\bar{x}| \times |\bar{y}| matrix of zeros.
 4:
 5:
           for (i, j) \in C do
                 if i = |\bar{x}| + 1 \lor j = |\bar{y}| + 1 then
 6:
                       continue
 7:
                 end if
 8:
                 if (rl_{\bar{x}}(i) = |\bar{x}| \wedge rl_{\bar{y}}(j) = |\bar{y}|) \vee c(x_i, y_j) = c(x_i, -) + c(-, y_j) then
 9:
                       if D_{i,j} = D_{i+1,j+1} + c(x_i, y_j) then
10:
                            \Gamma_{i,j} \leftarrow \Gamma_{i,j} + A_{i,j} \cdot B_{i+1,j+1}.
11:
                       end if
12:
                 else
13:
                      if D_{i,j} = D_{rl_{\bar{x}}(i)+1, rl_{\bar{y}}(j)+1} + d_{i,j} then
14:
                            \gamma \leftarrow \boldsymbol{A}_{i,j} \cdot \boldsymbol{B}_{rl_{\bar{x}}(i)+1,rl_{\bar{y}}(j)+1}.
15:
                            Compute D' and d' via Algorithm 5 for the subtrees \bar{x}^i and \bar{y}^j.
16:
                            D'_{1,2} \leftarrow \infty. \ D'_{2,1} \leftarrow \infty.
17:
                            (\mathbf{\Gamma}', k) \leftarrow \text{FORWARD-BACKWARD}(\bar{x}^i, \bar{y}^j, \mathbf{D}', \mathbf{d}', c).
18:
                            for i' \leftarrow 1, \ldots, |\bar{x}^i| do
19:
                                  for j' \leftarrow 1, \ldots, |\bar{y}^j| do
20:
                                       \Gamma_{i+i'-1,j+j'-1} \leftarrow \Gamma_{i+i'-1,j+j'-1} + \Gamma'_{i',j'} \cdot \gamma.
21:
                                  end for
22:
                            end for
23:
                       end if
24:
                 end if
25:
           end for
26:
           return (\Gamma, A_{|\bar{x}|+1, |\bar{y}|+1}).
27:
28: end function
```

Interestingly, the matrix Γ has further helpful properties. By considering the sum over all columns and subtracting it from the total number of co-optimal mappings we obtain the number of co-optimal mappings in which a certain node in \bar{x} is deleted. In our example, a is deleted in 2 co-optimal mappings, b in 3 co-optimal mappings, c in 4 co-optimal mappings, d in 4 co-optimal mappings, and e in 5 co-optimal mappings. Conversely, by summing up over all rows and subtracting the result from the total number of co-optimal mappings we obtain the number of co-optimal mappings in which a certain node of \bar{y} is inserted. In this example, neither f nor g are inserted in any co-optimal mappings.

Another interesting property is that the matrix Γ represents the frequency of certain pairings of nodes in co-optimal mappings, if we divide all entries by the total number of co-optimal mappings. This version of the matrix also offers an alternative view on the tree edit distance itself.

Theorem 21. Let \bar{x} and \bar{y} be trees over some alphabet \mathcal{X} , and let c be a cost function over \mathcal{X} . Further, let Γ and k be the two outputs of Algorithm 11 for \bar{x} , \bar{y} , and c, and let $P_c(\bar{x}, \bar{y}) := \frac{1}{k} \cdot \Gamma$. Then, the following equation holds:

$$d_{c}(\bar{x}, \bar{y}) = \sum_{i=1}^{|\bar{x}|} \sum_{j=1}^{|\bar{y}|} \mathbf{P}_{c}(\bar{x}, \bar{y})_{i,j} \cdot c(x_{i}, y_{j})$$

$$+ \sum_{i=1}^{|\bar{x}|} p_{i}^{\text{del}} \cdot c(x_{i}, -) + \sum_{j=1}^{|\bar{y}|} p_{j}^{\text{ins}} \cdot c(-, y_{j})$$

$$p_{i}^{\text{del}} := 1 - \sum_{j=1}^{|\bar{y}|} \mathbf{P}_{c}(\bar{x}, \bar{y})_{i,j}$$

$$p_{j}^{\text{ins}} := 1 - \sum_{i=1}^{|\bar{x}|} \mathbf{P}_{c}(\bar{x}, \bar{y})_{i,j}$$
(35)

Proof. Per construction, Γ is equivalent to the number of co-optimal mappings M, such that $(i,j) \in M$, and k is equivalent to the number of co-optimal mappings overall. The cost of each co-optimal mapping is per definition $d_c(\bar{x}, \bar{y})$. Therefore, summing over the cost of all these mappings and dividing by the number of mappings is also equal to $d_c(\bar{x}, \bar{y})$.

This alternative representation of the TED is particularly useful if one wishes to learn the parameters of the tree edit distance, as all the computational complexity of the tree edit distance is encapsulated in the matrix $P_c(\bar{x}, \bar{y})_{i,j}$ and all learned parameters are linearly multiplied with this matrix. This trick has been originally suggested by Bellet, Habrard, and Sebban (2012) to learn optimal parameters for the string edit distance.

The concludes our tutorial. For further reading, I recommend the robust tree edit distance by Pawlik and Augsten (2011), as well as the metric learning approaches by Bellet, Habrard, and Sebban (2012) and Paaßen et al. (2018).

References

Akutsu, Tatsuya (2010). "Tree Edit Distance Problems: Algorithms and Applications to Bioinformatics". In: *IEICE Transactions on Information and Systems* E93-D.2, pp. 208–218.

Bellet, Aurélien, Amaury Habrard, and Marc Sebban (2012). "Good edit similarity learning by loss minimization". In: *Machine Learning* 89.1, pp. 5–35. DOI: 10.1007/s10994-012-5293-8.

Choudhury, Rohan Roy, Hezheng Yin, and Armando Fox (2016). "Scale-Driven Automatic Hint Generation for Coding Style". In: *Proceedings of the 13th International Conference on Intelligent Tutoring Systems (ITS 2016)*. Ed. by Alessandro Micarelli, John Stamper, and Kitty Panourgia. Zagreb, Croatia, pp. 122–132. DOI: 10.1007/978-3-319-39583-8_12.

Freeman, Paul, Ian Watson, and Paul Denny (2016). "Inferring Student Coding Goals Using Abstract Syntax Trees". In: *Proceedings of the 24th International Conference on Case-Based Reasoning Research and Development (ICCBR 2016)*. Ed. by Ashok Goel, M Belén Díaz-Agudo, and Thomas Roth-Berghofer. Atlanta, GA, USA, pp. 139–153. DOI: 10.1007/978-3-319-47096-2_10.

- Henikoff, Steven and Jorja G. Henikoff (1992). "Amino acid substitution matrices from protein blocks".

 In: Proceedings of the National Academy of Sciences 89.22, pp. 10915–10919. URL: http://www.pnas.org/content/
- McKenna, Aaron et al. (2010). "The Genome Analysis Toolkit: A MapReduce framework for analyzing next-generation DNA sequencing data". In: *Genome Research* 20.9, pp. 1297–1303. DOI: 10.1101/gr.107524.110.
- Nguyen, Andy et al. (2014). "Codewebs: Scalable Homework Search for Massive Open Online Programming Courses". In: *Proceedings of the 23rd International Conference on World Wide Web (WWW 2014)*. Seoul, Korea, pp. 491–502. DOI: 10.1145/2566486.2568023.
- Paaßen, Benjamin (2018). "Metric Learning for Structured Data". under review. PhD thesis. Technical Faculty, Bielefeld University.
- Paaßen, Benjamin et al. (2018). "The Continuous Hint Factory Providing Hints in Vast and Sparsely Populated Edit Distance Spaces". In: Journal of Educational Datamining. accepted. URL: http://arxiv.org/abs/1
- Paaßen, Benjamin et al. (2018). "Tree Edit Distance Learning via Adaptive Symbol Embeddings". In: *Proceedings of the 35th International Conference on Machine Learning (ICML 2018)*. Ed. by Jennifer Dy and Andreas Krause. Vol. 80. accepted. Stockholm.
- Pawlik, Mateusz and Nikolaus Augsten (2011). "RTED: A Robust Algorithm for the Tree Edit Distance". In: *Proceedings of the VLDB Endowment* 5.4, pp. 334–345. DOI: 10.14778/2095686.2095692.
- (2016). "Tree edit distance: Robust and memory-efficient". In: *Information Systems* 56, pp. 157–173. DOI: 10.1016/j.is.2015.08.004.
- Rivers, Kelly and Kenneth R. Koedinger (2015). "Data-Driven Hint Generation in Vast Solution Spaces: A Self-Improving Python Programming Tutor". In: *International Journal of Artificial Intelligence in Education* 27.1, pp. 37–64. ISSN: 1560-4306. DOI: 10.1007/s40593-015-0070-z.
- Smith, Temple F. and Michael S. Waterman (1981). "Identification of common molecular subsequences". In: *Journal of Molecular Biology* 147.1, pp. 195–197. DOI: 10.1016/0022-2836(81)90087-5.
- Zhang, Kaizhong and Dennis Shasha (1989). "Simple Fast Algorithms for the Editing Distance between Trees and Related Problems". In: *SIAM Journal on Computing* 18.6, pp. 1245–1262. DOI: 10.1137/0218082.