

An A*-algorithm for the Unordered Tree Edit Distance with Custom Costs

Benjamin Paaßen^{*1}

¹Humboldt-University of Berlin

This is a preprint version of Paaßen (2021) as provided by the authors. For the original version, refer to Springer.

Abstract

The unordered tree edit distance is a natural metric to compute distances between trees without intrinsic child order, such as representations of chemical molecules. While the unordered tree edit distance is MAX SNP-hard in principle, it is feasible for small cases, e.g. via an A* algorithm. Unfortunately, current heuristics for the A* algorithm assume unit costs for deletions, insertions, and replacements, which limits our ability to inject domain knowledge. In this paper, we present three novel heuristics for the A* algorithm that work with custom cost functions. In experiments on two chemical data sets, we show that custom costs make the A* computation faster and improve the error of a 5-nearest neighbor regressor, predicting chemical properties. We also show that, on these data, polynomial edit distances can achieve similar results as the unordered tree edit distance.

Keywords: Unordered Tree Edit Distance; A* algorithm; Tree Edit Distance; Chemistry

1 Introduction

Tree structures occur whenever data follows a hierarchy or a branching pattern, like in chemical molecules (Gallicchio and Micheli, 2013; Rarey and Dixon, 1998), in RNA secondary structures (Shapiro and Zhang, 1990), or in computer programs (Paaßen et al., 2018). To perform similarity search on such data, we require a measure of distance over trees. A popular choice is the tree edit distance, which is defined as the cost of the cheapest sequence of deletions, insertions, and relabelings that transforms one tree to another (Bille, 2005; Zhang, 1996; Zhang and Shasha, 1989). Unfortunately, the edit distance becomes MAX SNP-hard for unordered trees, like tree representations of chemical molecules (Zhang and Jiang, 1994). Still, for smaller trees, we can compute the unordered tree edit distance (UTED) exactly using strategies like A* algorithms (Horesh et al., 2006; Yoshino et al., 2013). Roughly speaking, an A* algorithm starts with an empty edit sequence and then successively extends the edit distance such that a heuristic lower bound for the cost of the edit sequence remains as low as possible. The tighter our lower bound h , the more we can prune the search and the faster the A* algorithm becomes. Horesh et al. (2006) have proposed a heuristic based on the histogram

^{*}Funding by the German Research Foundation (DFG) under grant number PA 3460/2-1 is gratefully acknowledged.

of node degrees and Yoshino et al. (2013) have improved upon this heuristic by also considering label histograms and by re-using intermediate values via dynamic programming. However, both approaches assume unit costs, i.e. that deletions, insertions, and relabelings all have a cost of 1, irrespective of the content that gets deleted, inserted, or relabeled. This is unfortunate because, in many domains, we have prior knowledge that suggests different costs or we may wish to learn costs from data (Paaßen et al., 2018). Accordingly, most tree edit distance algorithms are general enough to support custom deletion, insertion, and replacement costs, as long as they conform to metric constraints (Bille, 2005; Zhang, 1996; Zhang and Shasha, 1989).

In this paper, we develop three novel heuristics for the A* algorithm which all support custom costs. The three heuristics have linear, quadratic, and cubic complexity, respectively, where the slower heuristics provide tighter lower bounds. Based on these novel heuristics, we investigate three research questions:

RQ1: Which of the novel heuristic is the fastest? And how do they compare against the state-of-the-art by Yoshino et al. (2013)?

RQ2: Do custom edit costs actually contribute to similarity search?

RQ3: How does UTED compare to polynomial edit distances in similarity search?

We investigate these research questions on two example data sets of chemical molecules, both represented as unordered trees. To answer RQ2 and RQ3, we consider a regression task where we try to predict the chemical properties of a molecule (boiling point and stability, respectively) via a nearest-neighbor regression. We begin our paper with more background and related work before we describe our proposed A* algorithm, present our experiments, and conclude.

2 Background and Related Work

Let Σ be an arbitrary set which we call *alphabet*. Then, we define a *tree* over Σ as an expression of the form $\hat{x} = x(\hat{y}_1, \dots, \hat{y}_K)$, where $x \in \Sigma$ and where $\hat{y}_1, \dots, \hat{y}_K$ is a list of trees over Σ , which we call the *children* of \hat{x} . If $K = 0$, we call $x()$ a *leaf*. We denote the set of all trees over Σ as $\mathcal{T}(\Sigma)$.

In this paper, we are concerned with similarity search on trees. In the literature, there are three general strategies to compute similarities on trees. First, we can construct a feature mapping $\phi : \mathcal{T}(\Sigma) \rightarrow \mathbb{R}^n$, which maps an input tree to a feature vector, and then compute a (dis-)similarity between features, e.g. via $d(\hat{x}, \hat{y}) = \|\phi(\hat{x}) - \phi(\hat{y})\|$. For example, we can represent trees by *pq*-grams (Augsten et al., 2008), by counts of typical tree patterns (Collins and Duffy, 2002), or by training a neural network (Gallicchio and Micheli, 2013; Kusner et al., 2017). The second strategy are tree kernels k , i.e. functions that directly compute inner products $k(\hat{x}, \hat{y}) = \phi(\hat{x})^T \cdot \phi(\hat{y})$ without the need to explicitly compute ϕ (Aiolli et al., 2015; Collins and Duffy, 2002).

In this paper, we focus on a third option, namely tree edit distances (Bille, 2005). Let Σ be an alphabet with $- \notin \Sigma$. Roughly speaking, a tree edit distance $d(\hat{x}, \hat{y})$ between two trees \hat{x} and \hat{y} from $\mathcal{T}(\Sigma)$ is the cost of the cheapest sequence of deletions, insertions, and relabelings of nodes in \hat{x} such that we obtain \hat{y} (Bille, 2005; Zhang, 1996; Zhang and Shasha, 1989). More precisely, let x_1, \dots, x_m be the nodes of \hat{x}^1 and y_1, \dots, y_n be the nodes of \hat{y}

¹Note that we use 'node' and 'label' interchangeably in this paper. To disambiguate between two nodes with the same label, we use the index i .

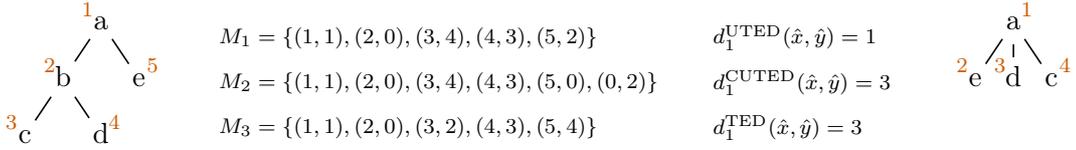


Figure 1: An illustration of mappings according to the unordered tree edit distance (Zhang and Jiang, 1994) (top), the constrained unordered tree edit distance (Zhang, 1996) (center), and the ordered tree edit distance (Zhang and Shasha, 1989) (bottom) between the same two trees. The distances assume unit costs. Numbers in superscript show the depth first order.

in depth-first-search order. Then, we define a *mapping* between \hat{x} and \hat{y} as a set of tuples $M \subset \{0, 1, \dots, m\} \times \{0, 1, \dots, n\}$ such that each $i \in \{1, \dots, m\}$ occurs exactly once on the left and each $j \in \{1, \dots, n\}$ occurs exactly once on the right. Figure 1 illustrates three example mappings between the trees $a(b(c, d), e)$ (left) and $a(e, d, c)$ (right), namely M_1 , M_2 , and M_3 (center left). Each mapping M can be translated into a sequence of edits by deleting all nodes x_i where $(i, 0) \in M$, by replacing nodes x_i with y_j where $(i, j) \in M$ and $x_i \neq y_j$, and by inserting all nodes y_j where $(0, j) \in M$. We denote the set of all mappings between \hat{x} and \hat{y} as $\mathcal{M}(\hat{x}, \hat{y})$. Next, we define a *cost function* as a metric $c : (\Sigma \cup \{-\}) \times (\Sigma \cup \{-\}) \rightarrow \mathbb{R}$, and we define the cost of a mapping M as $c(M) = \sum_{(i,j) \in M} c(x_i, y_j)$ where $x_0 = y_0 = -$. A typical cost function is $c_1(x, y) = 1$ if $x \neq y$ and $c_1(x, y) = 0$ if $x = y$, which we call *unit costs*. Finally, we define the tree edit distance $d_c : \mathcal{T}(\Sigma) \times \mathcal{T}(\Sigma) \rightarrow \mathbb{R}$ according to cost function c as the minimum $d_c(\hat{x}, \hat{y}) = \min_{M \in \mathcal{M}(\hat{x}, \hat{y})} c(M)$.

We obtain different edit distances depending on the additional restrictions we apply on the set of mappings $\mathcal{M}(\hat{x}, \hat{y})$. The unordered tree edit distance (UTED) requires that mappings respect the ancestral ordering, i.e. if $(i, j) \in M$, then descendants of i can only be mapped to descendants of j (Bille, 2005). A cheapest example mapping according to unit costs M_1 (Figure 1, top). The constrained unordered tree edit distance (CUTED) (Zhang, 1996) additionally requires that a deletion/insertion of a node implies either deleting/inserting all of its siblings or all of its children but one. This forbids M_1 and M_3 , where b is deleted but both its sibling and more than one child are maintained. M_2 is a cheapest mapping according to CUTED with unit costs. The ordered tree edit distance (TED) (Zhang and Shasha, 1989) requires that the ancestral ordering and the depth-first ordering is maintained. Accordingly, neither M_1 nor M_2 are permitted because they swap the order of c and d around. M_3 is a cheapest mapping according to TED with unit costs. Note that UTED is MAX-SNP hard. However, CUTED and TED are both polynomial (Zhang, 1996; Zhang and Shasha, 1989) via dynamic programming and we consider them as baselines in our experiments.

3 Method

In this section, we explain our proposed A* algorithm for the unordered tree edit distance (UTED). We begin with the general scheme, which we adapt from Yoshino et al. (2013), and then introduce three heuristics to plug into the A* algorithm.

A* algorithm: We first introduce a few auxiliary concepts that we require for the A* algorithm. First, let M be some subset of $\{0, \dots, m\} \times \{0, \dots, n\}$. Then, we denote with I_M the set $\{i > 0 \mid \exists j : (i, j) \in M\}$ and with J_M the set $\{j > 0 \mid \exists i : (i, j) \in M\}$, i.e. the set of left-hand-side and right-hand-side indices of M . Next, let \hat{x} and \hat{y} be trees with nodes x_1, \dots, x_m and y_1, \dots, y_n , respectively. Then, we define \mathcal{X}_i and \mathcal{Y}_j as the index sets of all descendants

Algorithm 1 The A* algorithm to compute the unordered tree edit distance $d_c^{\text{UTED}}(\hat{x}, \hat{y})$ between two trees \hat{x} and \hat{y} , depending on a cost function c and a heuristic h .

```

1: function ASTAR_UTED(trees  $\hat{x}$  and  $\hat{y}$ , cost function  $c$ , heuristic  $h$ )
2:   Initialize a priority queue  $Q$  with the partial mapping  $M = \{(1, 1)\}$ 
3:   and value  $c(x_1, y_1) + h(\{2, \dots, m\}, \{2, \dots, n\})$ .
4:   while  $Q$  is not empty do
5:     Poll partial mapping  $M$  with lowest value  $f$  from  $Q$ .
6:      $i \leftarrow \min\{1, \dots, m + 1\} \setminus I_M$ .
7:     if  $i = m + 1$  then
8:       return  $c(M \cup \{(0, j) \mid 1 \leq j \leq n, j \notin J_M\})$ .
9:     end if
10:    Retrieve  $(k, l) \in M$  with largest  $k$  such that  $x_k$  is ancestor of  $x_i$  and  $l > 0$ .
11:     $h^p \leftarrow h(\{1, \dots, m\} \setminus (\mathcal{X}_k \cup I_M), \{1, \dots, m\} \setminus (\mathcal{Y}_l \cup J_M))$ .
12:     $M_0 \leftarrow M \cup \{(i, 0)\}$ 
13:     $h_0 \leftarrow h(\mathcal{X}_k \setminus I_{M_0}, \mathcal{Y}_l \setminus J_{M_0}) + h^p$ .
14:    for  $j \in \mathcal{Y}_l \setminus J_M$  do
15:      Let  $y'_0, \dots, y'_t$  be the path from  $y_l$  to  $y_j$  in  $\hat{y}$  with  $y'_0 = y_l$  and  $y'_t = y_j$ .
16:       $M_j \leftarrow M \cup \{(i, j), (0, y'_1), \dots, (0, y'_{t-1})\}$ .
17:       $h_j \leftarrow h(\mathcal{X}_i \setminus I_{M_j}, \mathcal{Y}_j \setminus J_{M_j}) + h(\mathcal{X}_k \setminus (\mathcal{X}_i \cup I_{M_j}), \mathcal{Y}_l \setminus (\mathcal{Y}_j \cup J_{M_j})) + h^p$ .
18:    end for
19:    Put  $M_j$  with value  $c(M_j) + h_j$  onto  $Q$  for all  $j \in \{0\} \cup (\mathcal{Y}_l \setminus J_M)$ .
20:  end while
21: end function

```

of x_i and y_j , respectively. Finally, let c be a cost function. Then, we define a *heuristic* as a function $h : \mathcal{P}(\{1, \dots, m\}) \times \mathcal{P}(\{1, \dots, n\}) \rightarrow \mathbb{R}$, such that for any $I \subseteq \{1, \dots, m\}$ and $J \subseteq \{1, \dots, n\}$ it holds

$$h(I, J) \leq \min_{M \in \mathcal{M}^{\text{UTED}}(\hat{x}, \hat{y})} \sum_{(i, j) \in M: i \in I, j \in J} c(x_i, y_j). \quad (1)$$

Algorithm 1 shows the pseudocode for the A* algorithm. We initialize a partial mapping $M = \{(1, 1)\}$ which maps the root of \hat{x} to the root of \hat{y} . If this is undesired, input trees must first be augmented with a placeholder root node. Next, we initialize a priority queue Q with M and its lower bound. Now, we enter the main loop. In each step, we consider the current partial mapping M with the lowest lower bound f (line 5). If I_M already covers all nodes in \hat{x} , we complete M by inserting all remaining nodes of \hat{y} and return the cost of the resulting mapping (lines 7-9)². Otherwise, we extend M by mapping the smallest non-mapped index i either to zero (lines 12-13), or to j for some available node y_j from \hat{y} (lines 14-18). In the latter case, we need to maintain the ancestral ordering of the tree \hat{y} . Accordingly, we first retrieve the lowest ancestor x_k of x_i such that $(k, l) \in M$ and only permit i to be mapped to descendants \mathcal{Y}_l . Note that (k, l) must exist because we initialized M with $\{(1, 1)\}$. Further, if we map i to a non-direct descendant of y_l , we make sure to insert all nodes on the ancestral path y'_0, \dots, y'_t , first. We generate lower bounds h_j for all extensions M_j and put them back onto the priority queue.

Note that the space complexity of this algorithm can be polynomially limited by representing the partial mappings in a tree structure. However, the worst-case time complexity

²Strictly speaking, this is only valid if the lower bound f is exact for insertions. This is the case for all heuristics considered in this paper.

remains exponential because the algorithm may need to explore combinatorially many possible mappings. Generally, though, the tighter the lower bound provided by h , the fewer partial mappings need to be explored before we find a complete mapping. To further cut down the time complexity, we tabulate the lower bounds h^p for the ancestor mappings (k, l) (line 11), as recommended by Yoshino et al. (2013).

Heuristics: The final ingredient we need is the actual heuristic h . We define three heuristics in increasing relaxation and decreasing time complexity. First,

$$h_3(I, J) = \min_{M \subseteq \mathcal{M}(I, J)} \sum_{(i, j) \in M} c(x_i, y_j), \quad (2)$$

where $\mathcal{M}(I, J)$ denotes the set of all mappings between I and J , irrespective of ancestral ordering. Accordingly, Inequality 1 is trivially fulfilled because any mapping that respects ancestral ordering is also in $\mathcal{M}(I, J)$. Importantly, this relaxation can be solved in $\mathcal{O}((m+n)^3)$ via the Hungarian algorithm (Bougleux et al., 2017). While polynomial, this appears rather expensive for a heuristic. Therefore, we also consider further relaxations. Without loss of generality, let $|I| \geq |J|$, otherwise exchange the roles of \hat{x} , \hat{y} , I and J . Then, we define

$$h_2(I, J) = \min_{I' \subseteq I: |I'| = |I| - |J|} \left(\sum_{i \in I'} c(x_i, -) \right) + \left(\sum_{i \in I \setminus I'} \min_{j \in J} c(x_i, y_j) \right). \quad (3)$$

Note that this is a lower bound for $h_3(I, J)$ because we expand the class of permitted M to one-to-many mappings, which is a proper superset of $\mathcal{M}(I, J)$. Further, h_2 can be solved in $\mathcal{O}(m \cdot n)$ because we can evaluate $c_i := \min_{j \in J} c(x_i, y_j)$ for all i in $|I| \cdot |J|$ steps and we can solve the outer minimization by finding the $|I| - |J|$ smallest terms according to $c(x_i, -) - c_i$ and using those as I' , which is possible in $\mathcal{O}(|I|)$. In case even $\mathcal{O}(m \cdot n)$ is too expensive, we relax further to

$$h_1(I, J) = \min_{I' \subseteq I: |I'| = |I| - |J|} \sum_{i \in I'} c(x_i, -). \quad (4)$$

This is obviously a lower bound for h_2 and can be solved in $\mathcal{O}(\max\{m, n\})$.

4 Experiments

We evaluate our three research questions on two data sets from Chemistry, namely the Alkanes data set of 150 alkane molecules by Gallicchio and Micheli (2013) and the hundred smallest molecules from the ZINC molecule data set of Kusner et al. (2017). In the former case, the molecules are directly represented as trees (with 8.87 nodes on average) with hydrogen counts as node labels. In the latter case, we use the syntax tree of the molecule’s SMILES representation (with 13.82 nodes on average) (Weininger, 1988), where nodes are labeled with syntactic blocks. Note that this is a lossy representation because we cut aromatic rings to obtain trees.

Regarding RQ1, we compute all pairwise UTED values using the three heuristics h_1 , h_2 , and h_3 , both with unit costs and with custom costs. As custom cost function c , we use the difference in hydrogen count between two carbon for the alkanes data set. For the ZINC data set, we use the difference in electron count. For further reference, we also compare to the heuristic of Yoshino et al. (2013) for unit costs. We execute all computations in Python on a consumer desktop PC with Intel core i9-10900 CPU and 32 GB RAM and measure time using Python’s `time` function. All experimental code is available at <https://gitlab.com/bpaassen/uted>.

Table 1: The average runtime in milliseconds (top) and the number of partial mappings searched (bottom) per distance computation for each heuristic.

	data set	unit				custom		
		h_1	h_2	h_3	h_{yoshino}	h_1	h_2	h_3
runtime	alkanes	8.70	12.15	10.72	9.52	7.34	8.21	9.92
	ZINC	549.38	277.15	192.97	266.66	130.62	75.53	68.12
search size	alkanes	376	348	260	279	318	302	246
	ZINC	24586	9164	4158	6781	6643	2655	1379

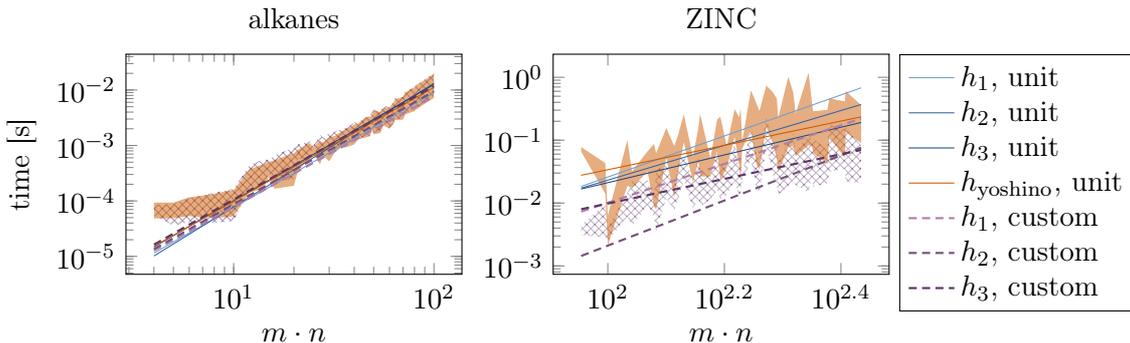
Figure 2: A log-log regression of the runtime needed for computing UTED for all four heuristics (indicated by color) on the alkanes data (left) and the ZINC data (right). Shading indicates distance between 25th and 75th percentile of the runtimes for h_{yoshino} (orange, solid), and h_3 with custom costs (purple, crosshatch), respectively.

Table 1 shows the average runtime in milliseconds (top) for each heuristic on both data sets. On alkanes, h_1 is fastest and on ZINC, h_3 is fastest. All heuristics get faster for custom costs. Surprisingly, h_{yoshino} is not the fastest for unit costs, even though it is optimized for this setting. This may just be due to an unfavourable constant factor, though: h_{yoshino} is successful in reducing the size of the search space almost to the same level as h_3 (see Table 1, bottom). Further, Figure 2 displays a linear regression for the runtime versus tree size in a log-log plot, indicating that h_{yoshino} and h_3 have the lowest slopes/best scaling behavior for large trees.

Regarding RQ2 and RQ3, we perform a 5-nearest neighbor regression³ to predict the boiling point of alkanes and the chemical stability measure of Kusner et al. (2017) for ZINC molecules, respectively. Table 2 shows the prediction error for both data sets in 15-fold cross-

³We also tested lower K , which achieved worse results for all methods.

Table 2: Average RMSE (\pm std.) of a 5-NN regressor across 15 crossvalidation folds for UTED, CUTED, and TED with unit and custom costs.

data set	unit			custom		
	UTED	CUTED	TED	UTED	CUTED	TED
alkanes	0.27 \pm 0.24	0.27 \pm 0.24	0.27 \pm 0.24	0.25 \pm 0.24	0.25 \pm 0.24	0.25 \pm 0.24
ZINC	1.33 \pm 0.85	1.31 \pm 0.86	1.36 \pm 0.84	1.24 \pm 0.87	1.26 \pm 0.87	1.29 \pm 0.86

validation. For reference, we do not only evaluate UTED with unit and custom costs, but also CUTED and TED. We observe that all methods perform better with custom costs compared to unit costs. For alkanes, there is no measurable difference between UTED, CUTED, and TED. For ZINC, TED performs worst and CUTED performs better than UTED for unit costs and UTED performs better than CUTED for custom costs.

5 Conclusion

We proposed three novel heuristics to compute the unordered tree edit distance via an A* algorithm. In contrast to prior work, our heuristics can accommodate custom costs, not only unit costs. Our three heuristics provide different trade-offs of time complexity (linear, quadratic, cubic) versus how much they prune the A* search.

In our experiments on two chemical experiments, we observed that this trade-off works in favor of the linear heuristic for small trees but that the cubic heuristic takes over for larger trees. Interestingly, the cubic heuristic compared favorably even to the current state-of-the-art heuristic. When applying custom costs, all our heuristics became faster thanks to the disambiguation provided by the custom cost function.

Regarding similarity search, we investigated the performance of a 5-nearest neighbor regressor, predicting chemical properties. We observed that custom costs lowered the regression error. However, we also saw that a similar performance can be achieved with a polynomial, restricted edit distance. Future work might investigate further tree data set to check whether these results generalize beyond chemistry.

References

- B. Paaßen, An A*-algorithm for the unordered tree edit distance with custom costs, in: N. Kriege, R. Connor, N. Reyes (Eds.), *Proceedings of the 14th International Conference on Similarity Search and Applications (SISAP 2021)*, Springer, 2021. Accepted.
- C. Gallicchio, A. Micheli, Tree echo state networks, *Neurocomputing* 101 (2013) 319 – 337. doi:10.1016/j.neucom.2012.08.017.
- M. Rarey, J. S. Dixon, Feature trees: a new molecular similarity measure based on tree matching, *Journal of computer-aided molecular design* 12 (1998) 471–490.
- B. A. Shapiro, K. Zhang, Comparing multiple RNA secondary structures using tree comparisons, *Bioinformatics* 6 (1990) 309–318. doi:10.1093/bioinformatics/6.4.309.
- B. Paaßen, C. Gallicchio, A. Micheli, B. Hammer, Tree Edit Distance Learning via Adaptive Symbol Embeddings, in: J. Dy, A. Krause (Eds.), *Proceedings of the 35th International Conference on Machine Learning (ICML 2018)*, volume 80 of *Proceedings of Machine Learning Research*, 2018, pp. 3973–3982. URL: <http://proceedings.mlr.press/v80/paassen18a.html>.
- P. Bille, A survey on tree edit distance and related problems, *Theoretical Computer Science* 337 (2005) 217 – 239. doi:10.1016/j.tcs.2004.12.030.
- K. Zhang, A constrained edit distance between unordered labeled trees, *Algorithmica* 15 (1996) 205–222. doi:10.1007/BF01975866.
- K. Zhang, D. Shasha, Simple fast algorithms for the editing distance between trees and related problems, *SIAM Journal on Computing* 18 (1989) 1245–1262. doi:10.1137/0218082.

- K. Zhang, T. Jiang, Some max snp-hard results concerning unordered labeled trees, *Information Processing Letters* 49 (1994) 249–254. doi:10.1016/0020-0190(94)90062-0.
- Y. Horesh, R. Mehr, R. Unger, Designing an A* algorithm for calculating edit distance between rooted-unordered trees, *Journal of Computational Biology* 13 (2006) 1165–1176. doi:10.1089/cmb.2006.13.1165.
- T. Yoshino, S. Higuchi, K. Hirata, A dynamic programming A* algorithm for computing unordered tree edit distance, in: 2013 Second IIAI International Conference on Advanced Applied Informatics, 2013, pp. 135–140. doi:10.1109/IIAI-AAI.2013.71.
- N. Augsten, M. Böhlen, J. Gamper, The pq-gram distance between ordered labeled trees, *ACM Transactions on Database Systems* 35 (2008). doi:10.1145/1670243.1670247.
- M. Collins, N. Duffy, New ranking algorithms for parsing and tagging: Kernels over discrete structures, and the voted perceptron, in: P. Isabelle, E. Charniak, D. Lin (Eds.), *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics (ACL 2002)*, 2002, pp. 263–270. URL: <http://www.aclweb.org/anthology/P02-1034.pdf>.
- M. J. Kusner, B. Paige, J. M. Hernández-Lobato, Grammar variational autoencoder, in: D. Precup, Y. W. Teh (Eds.), *Proceedings of the 34th International Conference on Machine Learning (ICML 2017)*, 2017, pp. 1945–1954. URL: <http://proceedings.mlr.press/v70/kusner17a.html>.
- F. Aioli, G. Da San Martino, A. Sperduti, An efficient topological distance-based tree kernel, *IEEE Transactions on Neural Networks and Learning Systems* 26 (2015) 1115–1120. doi:10.1109/TNNLS.2014.2329331.
- S. Bougleux, L. Brun, V. Carletti, P. Foggia, B. Gaüzère, M. Vento, Graph edit distance as a quadratic assignment problem, *Pattern Recognition Letters* 87 (2017) 38–46. doi:10.1016/j.patrec.2016.10.001.
- D. Weininger, Smiles, a chemical language and information system. 1. introduction to methodology and encoding rules, *Journal of Chemical Information and Computer Sciences* 28 (1988) 31–36. doi:10.1021/ci00057a005.