

APPROXIMATE NEAREST NEIGHBOR SEARCH FOR LABELLED TREES

LÁSZLÓ KOVÁCS Department of Information Technology, University of Miskolc kovacs@iit.uni-miskolc.hu

TIBOR RÉPÁSI Department of Information Technology, University of Miskolc repasi@iit.uni-miskolc.hu

ERIKA BAKSA-VARGA Department of Information Technology, University of Miskolc iitev@uni-miskolc.hu

[Received October 2004 and accepted January 2005]

Abstract. In many scientific areas there is a frequent need to extract a common pattern from multiple data. In most cases, however, an approximate but low cost solution is preferred to a high cost exact match. To establish a fast search engine an efficient heuristic method should be implemented. Our investigation is devoted to the approximate nearest neighbor search (ANN) for unordered labeled trees. The proposed modified best-first algorithm provides a $O((N_q+N_b)\cdot M + K\cdot N_q\cdot N_b/M)$ cost function with simple implementation details. According to our test results, realized with smaller trees where the brute-force algorithm could be tested, the yielded results are a good approximation of the global optimum values.

Keywords: tree matching, approximate nearest neighbor search

1. INTRODUCTION

In many scientific areas there is a frequent need to extract a common pattern from multiple data. The most common structure of the data is the hierarchy or tree. The task is to determine the set of sub-trees having the best matching with the pattern. Some important application areas for sub-tree matching are

pattern recognition, where the objects (for example pictures) are described by a tree structure.

molecular biology, where the real topology of RNA and of other molecules is a tree. From the topological similarities it is often possible to infer similarities in the function of the molecules. natural language processing, computational linguistics where dictionary definitions are stored in a lexical database. The definitions are represented syntactically as trees.

programming languages, where one of the main metadata structures is the tree structure. The parse trees, the operation trees or syntax trees are often used in the algorithms.

information systems, where the most common new storage format is the XML tree. The XML is seeing increased use and promises to fuel even more applications in the future. An XML document can be modeled as a tree. Each node in this tree corresponds to an element in the document. Each edge represents inclusion of the element corresponding to the child node under the element corresponding to the parent node in the XML file.

In the applications mentioned above the nodes of a tree are characterized by one or more attributes. Such attributes can be the type of the molecule (a node corresponds to a molecule) or the type of the operation (a node corresponds to an operation). The description vector of the nodes is called the label of the nodes. In some areas, not only the node types but also the order of nodes is important. In this case, the children are assigned to an ordering number in the scope of the parent. XML documents, for instance have an ordered and labeled tree structure. In our investigation we are focusing on unordered labeled tree structures. A recent workshop report from Yale suggested that more research should be undertaken to improve the heuristic search using algorithms designed to meet the demand made by increasingly large tree datasets [1].

2. RELATED WORKS

In this section we review the different researched approaches to comparing trees, as well as the algorithms developed so far to solve these problems. **P. Bille** published an extensive survey [2] on comparing trees with exact searching methods. As a conclusion from his work it turned out that all of the unordered versions of the problems in general are NP-hard. Indeed, the tree edit distance and alignment distance problems are even MAX SNP-hard. However, using special constraints polynomial time algorithms are available, just like for the ordered versions of the problems. These are all based on the classic technique of dynamic programming.

The general ordered tree edit distance, also called tree-to-tree correction problem was also fully reviewed in Technical Report 95-372 [3]. The problem was introduced by **Tai** [4] as a generalization of the string edit distance problem. His algorithm, which solved the problem without recursion, has its time and space complexity in $O(nmd^2d'^2)$, where n and m are the maximum number of children from any node in each of the trees, while d and d' are the maximum depth of the

trees. **Zhang and Shasha** [5] numbered the trees using postorder traversal instead of preorder, so the algorithm's space complexity is O(nm), while its time complexity is O(nmdd'). **Klein** [6] solved the problem in $O(n^3logn)$ time and O(nm) space. In his paper he proved that the algorithm can be extended to unrooted ordered trees within the same time and space bounds. **Chen** [7] applied fast matrix multiplication to solve the problem. The *unordered* version of the problem is NPcomplete even for binary trees with a label alphabet of size 2. It was shown in [8] that under special restrictions polynomial time algorithms exist.

There are other variants of the edit distance problem as well. One of them is the *unit cost edit distance*, where unit cost is defined as the number of edit operations required. In [9] the ordered version of the problem is considered and an algorithm with $O(u^2min\{n,m\}min\{l,l'\})$ time need is introduced, where *l* and *l'* are the number of leaves of the trees. The algorithm uses techniques from Ukkonen [10], and Landau and Vishkin [11]. The recursive solution of Selkow [12] used the basic operations, but insertions and deletions were restricted to leaf nodes only, which made the algorithm very simple and therefore its time complexity is O(nm). This is therefore sometimes referred to as the *1-degree edit distance*. Chawathe [13] utilizes the same restrictions, but in cases when external memory is needed to calculate the edit distance.

Tree inclusion, a special case of edit distance, is the problem to decide if tree T_1 can be included in T_2 . T_1 is included in T_2 if there is a sequence of delete operations performed on T_2 which make T_2 isomorphic to T_1 . For the ordered tree inclusion problem **Kilpeläinen and Mannila** [14] presented the first polynomial time algorithm using O(nm) time and space. A more space efficient version of this was given in [15] using O(nd') space. Later **Richter** [16] and **Chen** [17] developed more complex algorithms. In [14], [18] it is shown that the unordered tree inclusion problem is NP-complete. In spite of this an algorithm using $O(mni2^{2i})$ time exists.

Torsello and Hancock [19] prove, that a tree t' can be generated from a tree t with a sequence of node removal operations if and only if t' is an obtainable subtree of the directed association graph. Consequently the minimum cost edited tree isomorphism between two trees is a maximum common consistent subtree of the two directed association graphs if the node removal cost is uniform, and this result can also be extended to non-uniform cost. The background for this lies in [20], where the relationship between graph edit distance and the size of the maximum common subgraph is shown, and also their computational equivalence is demonstrated. This is an important observation since it has been established by **Barrow** and **Burstall** [21] that the maximum common subgraph problem may be transformed into a maximum clique problem using a derived structure referred to as the association graph. Pelillo et al. [22], for instance, transform the tree

isomorphism problem into a single max clique problem, a technique already used for the generic graph isomorphism problem. To obtain a maximal tree match, i.e. a maximal solution to the max clique problem, they use relaxation labeling. **Wang et al.** [23] considers the *largest approximately common subtree* problem for ordered, labeled trees using the edit distance to measure the dissimilarity of two trees. They present a dynamic programming algorithm, which runs as fast as the fastest known algorithm for computing the edit distance of trees.

This problem was investigated for unordered trees by Khanna, Motwani and Yao [24]. They created an algorithm for trees of bounded degree with performance ratio $O(nloglogn/log^2n)$ and then extended this to trees of unbounded degree with at most poly-log labels, obtaining a ratio of $O(n(loglogn)^2/log^2n)$. Akutsu and Halldórsson [25] also considers the approximation of the largest common subtree (and its special variation, the largest common edge subgraph) and largest common point set problems for unordered trees (and for ordered trees as a special case), and a general search algorithm is presented which approximates both problems within a factor of O(n/logn). For trees of bounded degree an improved algorithm is developed which approximates the largest common subtree within a factor of $O(n/log^2n)$. A large amount of work has been performed for comparing unordered trees based on various distance measures, especially on edit distance as the most commonly used distance measure. Shasha et al. [26], however, proposed a new approach, called Atree-Grep. They addressed the approximate nearest neighbor search problem for unordered labeled trees. Their algorithm, called 'pathfix', consists of two phases. First, the paths of the trees are stored in a suffix array and then the number of mismatching paths are counted between the query tree and the data tree. To speedup the search, they use a hash-based technique to filter out unqualified data trees at an early stage of the search. The algorithm has been implemented into two special Web-based search engines and proved to be fast, particularly when the dictionary size of node labels is large.

Other widely researched problems include tree pattern matching [27, 28, 29, 30, 31], maximum agreement subtree [32, 33] and smallest common supertree [34, 35].

3. DISTANCE MEASURES FOR TREE COMPARISON

As can be seen, most of the proposals in subtree matching are based on the edit distance between trees. This distance metric is a natural extension of the edit distance concept used for string comparisons. This metric provides an exact distance measurement between the trees. The drawback of these algorithms is the high cost of the computations. In the case of online applications with large tree datasets, the execution time is a crucial factor. In these kinds of applications, an approximate but low cost solution is preferred to a high cost exact solution. Our investigation is devoted to the approximate nearest neighbor search (ANN) for unordered labeled trees. Our goal is to construct an efficient heuristic method for the ANN problem. Since the ANN problem for edit distance metric is an NPproblem as is proven in [8], a modified distance definition is introduced.

Let D denote a domain set. This contains the node labels. The symbol T denotes an unordered, labeled tree. The following denotations related to the tree structure are used in the paper:

n	a node of the tree
l(n)	the label of node n , $l(n) \in D$
T_D	set of unordered, labeled trees on D
V(T)	vertices of T
E(T)	edges of T
r(T)	the root node of T

The goal of the investigation is to find the neighboring trees based on the similarity values. The distance or similarity is usually measured by a metric function. A space X is called metric space if a d(A,B) real non-negative function of two objects is defined with the following properties:

For every $A \in X$, $B \in X$, and $C \in X$.

- 1. d(A, B) = 0 if and only if A = B (the distance is 0 if and only if the points coincide),
- 2. d(A, B) = d(B, A) (the distance from A to B is the same as the distance from B to A),
- 3. $d(A,B) + d(B,C) \ge d(A,C)$ (the sum of two sides of a triangle is never less than the third side).

The d(A,B) function is known as the distance between the two points.

In the case of edit distance, a set of elementary transformation functions is defined on T_D . This set is denoted as E_D . The cost value of the elementary transformations is a non-negative real number. The corresponding cost function is denoted by

$$c \quad E_D \to R^+$$

It is assumed that T_D is closed to E_D , i.e.

$$e \quad T_D \to T_D \quad \forall e \in E_D,$$

$$\forall T_1, T_2 \in T_D \ \exists e_1, e_2, \dots, e_m \in E_D: e(T_1) = e_m \ o \ e_{m-1} \ o \dots e_2 \ o \ e_1(T_1) = T_2.$$

Let us denote the set of chain of transformations from T_i to T_j by $E_{i,j}$. The cost of chain e is defined as the sum of the single transformation steps:

$$c(e) = \Sigma c(e_i).$$

The edit distance between T_i and T_j is defined as the minimal cost of transformation chains from T_i to T_j :

$$c_{i,j} = \min\{c(e) \mid e \in E_{i,j}\}.$$

Usually, like in [36] the following elementary e operations are defined for tree objects:

relabel: assigns a new node name to the root of the tree, *insert*: inserting a new node into the children of the root node, *delete*: deleting a node from the children of the root node, *insert tree*: inserting a tree under the root node, *delete tree*: deleting a tree from the children of the node.

The list of elementary transformations with minimal cost is usually generated with a dynamic programming method. According to [2, pp.7], the tree distance value can be calculated using the following recursive formula:

$$d(0,0) = 0$$

$$d(F,0) = d(F-v,0) + c(v,0)$$

$$d(0,F) = d(0,F-v) + c(0,v)$$

$$d(F_1,F_2) = min \begin{cases} d(F_1-v,F_2) + c(T(v),0) \\ d(F_1,F_2-v) + c(0,T(v)) \\ d(F_1-T(v),F_2-T(w)) + c(T(v),T(w)) \end{cases}$$

where F denotes a tree, T(v) denotes a tree with root element v, and c(x,y) denotes the cost of transforming node x to node y. The computation cost of the basic dynamic programming method for trees is $O(|T|^4)$. This is a very high cost value for an ANN problem, as the distance computation should be calculated for a large number of pairs. It is proved in [26] that the ANN problem for edit distance metric is an NP-complete problem. In spite of this difficulty, most of the proposals for ANN searching for trees use the edit distance measure. There are very few proposals that apply a simplified distance function to provide a lower cost solution.

A good example for this approach is [26], where the distance from T_1 to T_2 is measured with the total number of root-to-leaf paths in T_1 that do not appear in T_2 . The nodes in T_2 that do not appear in T_1 can be freely removed. As can be seen, this definition introduces an asymmetric distance concept. In the definition T_1 denotes the query tree while T_2 is the searched tree. In our approach, another simplified distance function was selected.

4. MODIFIED BEST-FIRST ALGORITHM

Two trees are said to be similar if they have similar vertices with similar edges. During the editing process every vertex of the query tree is either transformed into a vertex of the base tree or it is deleted. Based on this transformation, every vertex of the query tree can be mapped either to a target vertex or to the sink symbol. Using this approach, a generalized mapping can be defined between the query and the base tree. We define m() as a monomorphism from T_1 to T_2 in the following way:

- 1. m $V(T_1) \rightarrow V(T_2) \cup \epsilon$
- 2. $\forall v, m(v) \in V(T_2)$ l(v) = l(m(v))
- 3. $\forall v_1 \neq v_2, m(v_1), m(v_2) \in V(T_2)$ $m(v_1) \neq m(v_2)$
- 4. $\forall v_1 \neq v_2, m(v_1), m(v_2) \in V(T_2): v_1 \leq v_2 \Leftrightarrow m(v_1) \leq m(v_2)$

According to the first property, every node in T_1 is mapped either to a node in T_2 or is deleted, i.e. it is mapped to the ε symbol. The second property says that a vertex should be mapped only to nodes of the same label. Due to the third property, the different query vertices can not be mapped to the same base vertex. The fourth property is called ancestor condition, the ancestor-descendants relationship among the query vertices must be preserved in the target tree, too.

Other types of relationships among the query vertices are neglected and not preserved. In this approach, the sibling vertices may be mapped to parent-child vertices, if the existing parent-child relationships are preserved. The parent-child relationships are the only important information stored in the query tree. The absence of an edge means in our approach a 'do not know' information. In this case, we don not care about the existence of an edge between the mapped vertices in the base tree. Figure 1 shows an example for this mapping.



Figure 1: Distance mapping example

Figure 1a) and Figure 1b) show valid mappings. The sibling nodes in the query tree are mapped to sibling nodes in Figure 1a), and to parent-child nodes in Figure 1b). Figure 1c) shows an invalid mapping as the parent-child relationship is not preserved. This kind of similarity value differs from the usual edit distance in the following aspects: 1) it does not take the re-labeling operation into account, and 2) only one side of the operands can be deleted.

Based on this mapping, a similarity value can be defined between two trees. The cost of mapping m is defined as the sum of the vertex mappings related to the query tree:

$$cost(m) = \sum_{n \in V(T)} c(n),$$

where

$$c(n) \begin{cases} C_2, \text{ if } m(n) = \varepsilon \lor m(r(T)) = \varepsilon \\ 0, \text{ if } n = r(T) \land m(r(T)) \neq \varepsilon \\ C_1 \cdot (d(m(n), m(a(n)) - 1) \text{ otherwise.} \end{cases}$$

In this definition, a(n) denotes the nearest ancestor of n in the query tree which is mapped to a non- ε element. If the root of the query tree is mapped to ε then c(n) is C_2 , otherwise the path from n to r(T) (excluding n and including r(T)) contains minimum one vertex mapped to a non- ε value. In this case both m(n) and m(a(n))are non- ε elements. The d() function denotes the length of path from m(a(n)) - m(n) in the base tree. As mapping m preserves the parent-child relationship, m(a(n)) is an ancestor of m(n). Thus d() yields a positive integer value. C_1 and C_2 are cost units. C_1 corresponds to gap-lengths between two preserved vertices and C_2 denotes the cost for vertex deletion. In our approach, C_2 is greater than C_1 since the absence of an element means a larger difference than the relocation of the element.

In this definition, a(n) denotes the nearest ancestor of n in the query tree which is mapped to a non- ε element. If the root of the query tree is mapped to ε then c(n) is C_2 , otherwise the path from n to r(T) (excluding n and including r(T)) contains minimum one vertex mapped to a non- ε value. In this case both m(n) and m(a(n))are non- ε elements. The d() function denotes the length of path from m(a(n)) - m(n) in the base tree. As mapping m preserves the parent-child relationship, m(a(n)) is an ancestor of m(n). Thus d() yields a positive integer value. C_1 and C_2 are cost units. C_1 corresponds to gap-lengths between two preserved vertices and C_2 denotes the cost for vertex deletion. In our approach, C_2 is greater than C_1 since the absence of an element means a larger difference than the relocation of the element.



Figure 2: An example for mapping

As an example, let the calculation of the mapping cost for Figure 2 stay here. The cost for root mapping is 0. The cost for white node is also 0 (there is no gap in the

mapped path). The cost for black node is 1 (one node length gap). The total cost is 0 + 0 + 1 = 1. We remark that in some applications it seems useful to introduce a weight factor in the cost expression. In this case the different edges may have different importance factors.

It can be seen that the distance measure based on this cost value does not meet the requirements of a metric space. The metric distance function should be symmetric while the given cost function is asymmetric. The roles of the query and base trees are distinguished. This corresponds to our intention, as we try to find a best matching sub-tree included in the base tree. The goal is to find a mapping with minimal cost value.

Taking a query tree T_l with N_q nodes and a base tree T_2 with N_b nodes, the number of potential mappings is O(N_b! / (N_b-N_q)!). Although the ancestor criteria restricts the set of potential mappings, the number of possible enabled mappings is too high. It would be very costly to test all of the possible mappings. Thus some kind of heuristics should be applied to speed up the matching process. In our investigation a variant of the best-first search method was selected.

The best first search method works on a state-tree. Each node of the tree is assigned to a cost value. The goal is to find the path with the minimal cost value. The best-first search divides the nodes into three distinct groups: the nodes tested (G_1) , the nodes ready to be tested (G_2) , and the rest (G_3) . Initially, G_1 is empty and G_2 contains only the root element. In a loop, the node from the ready state with the best (minimal) cost value is selected to be tested. During the test, the children of the node are evaluated and moved from the G_3 group into the G_2 group. The loop terminates if a leaf node is selected for testing.

In the applied variant, the nodes of the state-tree are assigned not to the vertices but to the vertex mappings of the query tree. Thus each node represents a decision about the mapping function. The state-tree is expanded and traversed in the following way:

- 1. Generating a label vector for every node. The label vector contains the counter values for the different labels related to the nodes in the descendant set. This vector works similar to one-grams used in the string distance problem. In the example shown in Figure 3a), the description vector for the root node is lv(3,2,1,3), where the first dimension is assigned to the green label, the second to the red label, the third to the blue label and the fourth to the black label.
- 2. Calculation of the label vectors for the query tree.
- 3. Selecting maximum K nodes in the base tree with the same label as the root of the query tree and with the first K best similarity values regarding the label vectors. The similarity value for label vectors is defined by

$$d(l_q, l_b) = \sum_j \max(l_{qj} - l_{bj}, 0)$$

where l_q belongs to the query tree and l_b to the base tree.

- 4. Loop on the selected nodes. Let w denote the vertex actually tested. Map the root of the query tree to w. Empty G_2 and G_1 .
- 5. Insert the mapping of w into G_2 .
- 6. Take the element x from G_2 with the lowest cost value. Move x from G_2 into G_1 . Disable the other mappings in G_2 from x or to m(x).
- 7. Test the children vertices of x considering the query tree. For every vertex generate the set of possible mappings. Evaluate these mappings and insert them into G_2 .
- 8. If G_2 is empty, the procedure terminates. The sum of cost values for the selected mappings is the approximation of the best mapping cost value for w, denoted by C(w). Go back to step 4.
- 9. Return $min\{C(w)\}$ as the approximation of the optimal mapping cost.

The cost of generating the label vectors is $O((N_q+N_b)M)$ as every vertex should be accessed only once. The label vector of a node can be built from the label vectors of its children. In the cost expression M denotes the number of different label values. M corresponds to the length of the label vectors. During the best-first search N_q vertices are tested and expanded. A vertex from the query tree may be mapped to $O(N_b/M)$ vertices in average. As the best-first search is repeated by Ktimes, the cost estimation for the algorithm is $O((N_q+N_b)\cdot M + K\cdot N_q\cdot N_b/M)$. Thus the cost is linear in both N_q and N_b . This cost is a significant reduction compared with the $O(M\cdot N_b! / (N_b-N_q)!)$ value for the brute force search method.

5. RESULTS

The implementation tests show a similar linearity for the computation costs. The test programs are implemented in the Scilab language. The next small example illustrates the cost relations between the brute-force and the heuristic method. The base tree has 10 vertices and is shown in Figure 3a). The query tree has 4 vertices and is shown in Figure 3b). The number of labels is 4. The trees were generated randomly.



Figure 3: Mapping example

The elements of the best mapping are given in the Figure with blue arrows. The cost value is only 1. Both methods can detect this optimal mapping but with a very different cost value. Table 1 shows the execution cost values for the investigated methods related to this example query.

•						
	Method	Cost				
	Brute-force	69.48 sec				
	M. Best-first	00.08 sec				
	Selkow	02.46 sec				

Table	1:	Comparing	the costs	of	execution
-------	----	-----------	-----------	----	-----------

To test the cost values for examples of larger tree sizes, a test run was implemented with values $N_q = 12$, $N_b \in [15..500]$. The cost values are shown in Figure 4.



Figure 4: Cost values for larger trees

DIGITALIZÁLTA: MISKOLCI EGYETEM KÖNYVTÁR, LEVÉLTÁR, MÚZEUM

The x-axis denotes the N_b value, while the y-axis shows the computation cost (where the maximum value is 3.6 sec). The trees were generated randomly. In Figure 4 the linearity of the cost function is well demonstrated.

6. COMPARISON WITH EDITING DISTANCE

Most of the works related to the comparison of unordered trees are based on their editing distance or Levenshtein distance. The Levenshtein distance is defined as the smallest number of insertions, deletions, and substitutions required to change one string or tree into another [37].

Most of the related works in this field are based on **Selkow**'s algorithm introduced in [12] and summerized by **P. Bille** in [2]. Selkow's algorithm calculates the editing distance between two forests of ordered trees. The measured editing distance is very similar to the editing distance of strings. It is simplified to a one level comparison, so edition is necessary each time the root nodes are not identical. As an ordered tree is a special case of an unordered tree, Selkow's algorithm can be used for unordered trees as well. It is a basic algorithm which needs a tremendous computation power of $O(N^2 \cdot M^2)$, where N and M are the node numbers of the trees, to find the editing order of the least cost. In case of unordered trees the algorithm has to take each order of the tree in account, so the computation cost will grow by $O(2^{N} \cdot 2^{M})$ for the repetitions on each possible permutation of the trees. The basic operations Selkow's algorithm makes use of node delete, node insert and node substitution or relabeling. Each operation can have different cost functions. The algorithm itself is recursive very much like above mentioned in section 3. It is working in the following way:

- 1. searching for the roots of the source forests,
- 2. extracting each tree of the forests,
- 3.a computing the cost (by doing a recursion) of deleting the root-node for each tree of one forest,
- 3.b computing the cost (by doing a recursion) of inserting a new root-node for each tree of the other forest,
- 3.c computing the cost (by doing a recursion) of substituting the root node for each tree of both forests,
- 4. selecting the minimum of all costs and returning it to the upper level of the recursion.

The generalization of the algorithm to forests is inevitable due to the fact, that the algorithm is recursive and that deleting the root node of a tree will result in a forest. We have tested an implementation of Selkow's algorithm on small trees. Running our implementation of the algorithm with the trees shown on Figure 3 we got the editing distance of 10 units, considering 1 as the cost of each edit operation. The

time needed to complete the calculation in the same environment was 2.463 sec as is shown in Table 1. Running the algorithm with randomly generated trees will show a very noisy cost function, however, the limits should show the $O(N^2 \cdot M^2)$ characteristics.

7. CONCLUSION

The approximate sub-tree search for trees with edit distance metric is an NPcomplete problem. To establish a fast search engine an efficient heuristic method should be implemented. The proposed modified best-first method provides a $O((N_q+N_b)\cdot M + K\cdot N_q\cdot N_b/M)$ cost function with simple implementation details. According to our test results, realized with smaller trees where the brute-force algorithm could be tested, the yielded results are a good approximation of the global optimum values.

REFERENCES

- [1] CRACRAFT, J. DONOGHUE, M.: Assembling the tree of life: Research needs in phylogenetics and phyloinformatics. Report from NSF Workshop, Yale University, July 2000.
- [2] BILLE, P.: Tree Edit Distance, Alignment Distance and Inclusion. IT University of Copenhagen, Technical Report Series TR-2003-23, ISSN 1660-6100, March 2003.
- [3] BARNARD, CLARKE, DUNCAN: Tree-to-Tree Correction for Document Trees. Technical Report 95-372, Queen's University Canada, January 1995.
- [4] TAI: *The Tree-to-Tree Correction Problem*. Journal of the Association for Computing Machinery (JACM), 26:422-433, 1979.
- [5] ZHANG, SHASHA: Simple fast algorithms for the editing distance between trees and related problems. SIAM Journal of Computing, 18:1245-1262, 1989.
- [6] KLEIN: Computing the edit-distance between unrooted ordered trees. In Proceedings of the 6th annual European Symposium on Algorithms (ESA) 1998, pp. 91-102.
- [7] CHEN: New algorithm for ordered tree-to-tree correction problem. Journal of Algorithms, 40:135-158, 2001.
- [8] ZHANG, STATMAN, SHASHA: On the editing distance between unordered labeled trees. Information Processing Letters, 42:133-139, 1992.
- [9] SHASHA, ZHANG: Fast algorithms for the unit cost editing distance between trees. Journal of Algorithms, 11:581-621, 1990.
- [10] UKKONEN: Finding approximate patterns in strings. Journal of Algorithms, 6:132-137, 1985.

- [11] LANDAU, VISHKIN: Fast parallel and serial approximate string matching. Journal of Algorithms, 10:157-169, 1989.
- [12] SELKOW: The tree-to-tree editing problem. Information Processing Letters, 6(6):184-186, 1977.
- [13] CHAWATHE: Comparing hierarchical data in extended memory. In Proceedings of VLDB, 1999, pp. 90-101.
- [14] KILPELÄINEN, MANNILA: Ordered and unordered tree inclusion. SIAM Journal of Computing, 24:340-356, 1995.
- [15] KILPELÄINEN: Tree Matching Problems with Applications to Structured Text Databases. PhD Thesis, University of Helsinki, Department of Computer Science, 1992.
- [16] RICHTER: A new algorithm for the ordered tree inclusion problem. In Proceedings of the 8th Annual Symposium on Combinatorial pattern Matching (CPM), in Lecture Notes of Computer Science (LNCS), vol. 1264, pp 150-166, Springer 1997.
- [17] CHEN: More efficient algorithm for ordered tree inclusion. Journal of Algorithms, 26:370-385, 1998.
- [18] MATOUSEK, THOMAS: On the complexity of finding iso- and other morphisms for partial k-trees. Discrete Mathematics, 108:343-364, 1992.
- [19] TORSELLO, HANCOCK: Computing approximate tree edit distance using relaxation labeling. Pattern Recognition Letters 2003, PII: S0167-8655(02)00255-6, 2002.
- [20] BUNKE, KANDEL: Mean and maximum common subgraph of two graphs. Pattern Recognition Letters 21, 2000, pp. 163-168.
- [21] BARROW, BURSTALL: Subgraph isomorphism, matching relational structures and maximal cliques. Information Processing Letters 4, 1976, pp. 83-84.
- [22] PELLILO et al.: Matching hierarchical structures using association graphs. IEEE PAMI 21, 1999, pp. 1105-1120.
- [23] WANG et al.: An Algorithm for Finding the Largest Approximately Common Substructures of Two Trees.
- [24] KHANNA, MOTWANI, YAO: Approximation algorithms for the largest common subtree problem. Technical Report, Stanford University, 1995.
- [25] AKUTSU, HALLDÓRSSON: On the Approximation of Largest Common Subtrees and Largest Common Point Sets. Science Institute University of Iceland, October 1997.
- [26] SHASHA et al.: AtreeGrep Approximate Searching in Unordered Trees. In Proceedings of SSDBM 2002, Edinburgh, July 2002, pp. 89-98.

- [27] KOSARAJU: *Efficient tree pattern matching*. In Proceedings of the 30th IEEE Symposium on the Foundations of Computer Science (FOCS), 1989, pp. 178-183.
- [28] DUBINER, GALIL, MAGEN: Faster tree pattern matching. In Proceedings of the 31st IEEE Symposium on the Foundations of Computer Science (FOCS), 1990, pp. 145-150.
- [29] HOFFMANN, DONNELL: Pattern matching in trees. Journal of the Association for Computing Machinery (JACM), 29(1):68-95, 1982.
- [30] RAMESH, RAMAKRISHNAN: Nonlinear pattern matching in trees. Journal of the Association for Computing Machinery (JACM), 39(2):295-316, 1992.
- [31] ZHANG, SHASHA, WANG: Approximate tree matching in the presence of variable length don't cares. Journal of Algorithms, 16(1):33-66, 1994.
- [32] KESELMAN, AMIR: Maximum agreement subtree in a set of evolutionary trees metrics and efficient algorithms. In Proceedings of the 35th Annual Symposium on the Foundations of Computer Science (FOCS), 1994, pp. 758-769.
- [33] FARACH, THORUP: Fast comparison of evolutionary trees. In Proceedings of the 5th Annual ACM-SIAM Symposium on Discrete Algorithms, 1994, pp. 481-488.
- [34] NISHIMURA, RAGDE, THILIKOS: Finding smallest supertrees under minor containment. International Journal of the Foundations of Computer Science, 11(3):445-465, 2000.
- [35] GUPTA, NISHIMURA: Finding largest subtrees and smallest supertrees. Algorithmica, 21:183-210, 1998.
- [36] NIERMAN, JAGADISH: Evaluating Structural Similarity in XML Documents. University of Michigan, IIS-0002356.
- [37] LEVENSHTEIN, V. I.: Binary codes capable of correcting deletions, insertions, and reversals. Doklady Akademii Nauk SSSR, 163(4):845-848, 1965 (Russian). English translation in Soviet Physics Doklady, 10(8):707-710, 1966.