

Algorithms for the Edit Distance Between Trees

Shaofeng Jiang

Western University

sjian7@uwo.ca

December 5, 2017

Overview

- 1 Introduction
- 2 Notation
- 3 Recursive Solution
- 4 A Simple Algorithm
 - Key Idea
- 5 Algorithmic Improvement
 - Time Complexity
- 6 Zhang and Shasha's algorithm
 - Key Idea
 - Key Concept
 - Time Complexity
- 7 Klein's Algorithm
 - Definition
 - Observation
 - Key Idea
 - Key Concept
- Algorithm
- Time Complexity
- 8 Decomposition Strategies
- 9 Cover Strategies
- 10 A New Optimal Cover Strategy
 - Key Idea
 - Algorithm
 - Time Complexity
 - Performance
- 11 An Optimal Decomposition Algorithm
 - Key Concept
 - Key Idea
 - Algorithm
 - Time Complexity
- 12 Optimal Decomposition Strategy
- 13 Implement

Ordered Labelled Tree

An ordered labelled tree is a tree in which the nodes are labeled and the left-to-right order among siblings is significant.

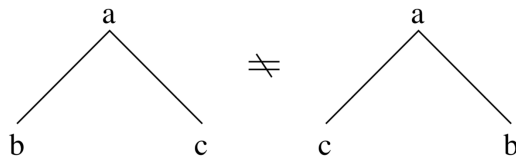


Figure: Ordered Labelled Tree

The ordered labelled tree can represent:

- an XML document
- a natural language parse
- RNA secondary structure
- and so on...

Pattern Matching in Trees

Given a Pattern tree and a Data tree, we may want to match the pattern tree to the data tree.

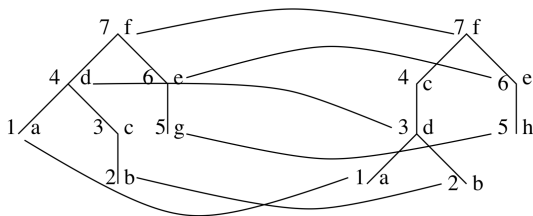


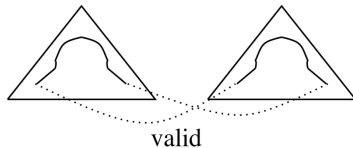
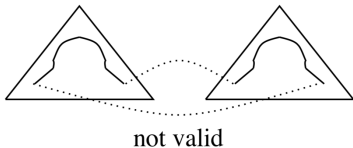
Figure: Pattern Matching in Trees

Pattern Matching in Trees

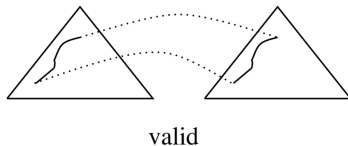
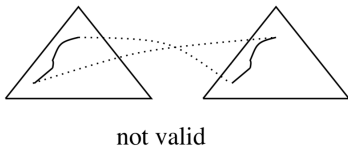
- **Input:** Two labelled ordered trees T and P .

- **Constraints:**

- One-to-one relationship
- Sibling order is preserved.



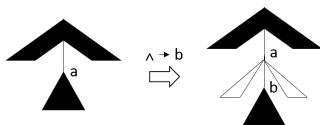
- Ancestor order is preserved.



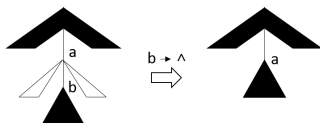
- **Output:** Mapping between T and P .

Edit Operations

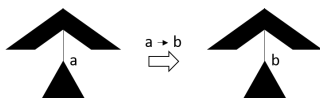
- Mappings on trees can be simplified into sequences of edit operations to transform one tree to another.
- Three valid edit operations supported:
 - **Insertion.** insertion of a node.



- **Deletion.** deletion of a node

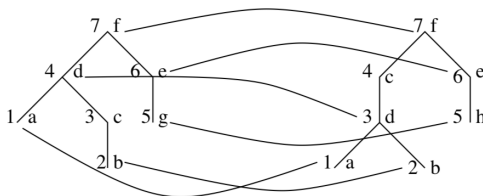


- **Substitution.** replace a label of a node by another label.



Cost of Mapping and the Edit Distance

- The cost of a mapping is the sum of its insertion, deletion and substitutions.



Edit sequence: $(c \rightarrow \lambda), (g \rightarrow h), (\lambda \rightarrow c)$

Figure: The Edit Sequence and its Cost

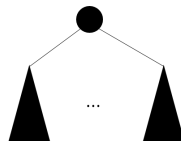
- The edit distance between two trees P and D is the minimum cost to change P to D via a sequence of basic edit operations.

- **Tree**

- $l(f)$



- $l(A_1 \circ A_2 \circ \dots \circ A_n)$



Let T be a tree,

- $r(T)$

the root of T

- T°

$$T - r(T)$$

- $lr(T)$

the root of the leftmost tree in the forest T°

- $rr(T)$

the root of the rightmost tree in the forest T°

- **Forest**

- $l(f) \circ t$

- $t \circ l(f)$



Let F be a forest,

- $|F|$ the number of nodes in the forest F
- $\#leaves(F)$ the number of leaves in the forest F
- $F(i)$ the sub-tree of F rooted at node i
- $F - i$ the forest after deleting node i
- $lr(F)$ the root of the leftmost tree in the forest F
- $rr(F)$ the root of the rightmost tree in the forest F

String Edit Distance and its Recursive Solution

- String edit distance is the minimum cost to change one string to another via a sequence of basic edit operations.
- The recursive solution of the string edit problem $d(S_1, S_2)$ is

$$d(S_1, S_2) = \min \begin{cases} d(S_1 - u, S_2) + \delta(u, \emptyset) \\ d(S_1, S_2 - v) + \delta(\emptyset, v) \\ d(S_1 - u, S_2 - v) + \delta(u, v) \end{cases}$$

String Edit Distance and its Recursive Solution

- If u and v is both the first element of string S_1 and S_2 , then it is a leftmost decomposition.

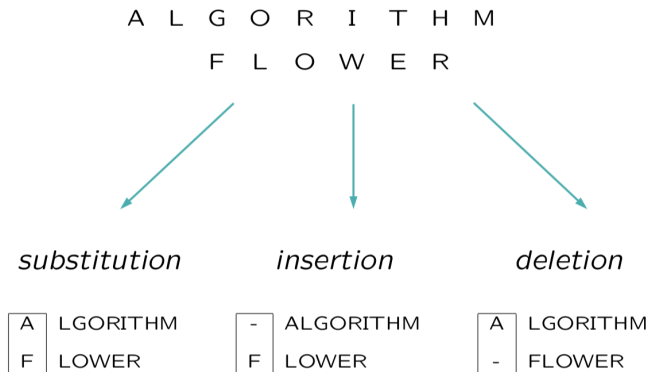


Figure: String Leftmost Decomposition

String Edit Distance and its Recursive Solution

- If u and v is both the last element of string S_1 and S_2 , then it is a rightmost decomposition.

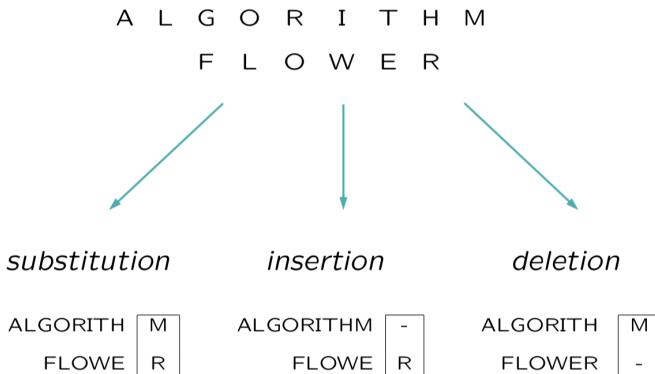


Figure: String Rightmost Decomposition

Tree Edit Distance and its Recursive Solution

For tree-to-tree distance, Let T_1 and T_2 are both trees of the form $l(f)$ and $l'(f')$,

$$d(l(f), l'(f')) = \min \begin{cases} d(f, l'(f')) + \delta(l, \emptyset) \\ d(l(f), l'(f')) + \delta(\emptyset, l') \\ d(f, f') + \delta(l, l') \end{cases}$$

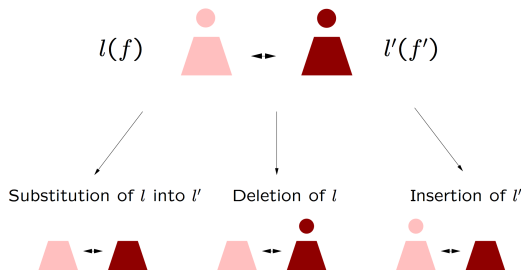


Figure: Tree-to-Tree Edit Distance

Tree Edit Distance and its Recursive Solution

For forest-to-forest distance, Let F_1 and F_2 are both trees of the form $l(f) \circ t$ and $l'(f') \circ t'$,

$$d(l(f) \circ t, l'(f') \circ t') = \min \begin{cases} d(f \circ t, l'(f') \circ t') + \delta(l, \emptyset) \\ d(l(f) \circ t, f' \circ t') + \delta(\emptyset, l') \\ d(f \circ t, f' \circ t') + \delta(l, l') \end{cases}$$



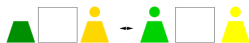
Leftmost decomposition

Substitution of l into l'



$$\begin{aligned} &\text{Distance}(l(f), l'(f')) \\ &+ \\ &\text{Distance}(t, t') \end{aligned}$$

Deletion of l



$$\text{Distance}(f \circ t, l'(f') \circ t') + \text{del}(l)$$

Insertion of l'



$$\text{Distance}(l(f) \circ t, f' \circ t') + \text{ins}(l')$$

Tree Edit Distance and its Recursive Solution

For forest-to-forest distance, Let F_1 and F_2 are both trees of the form $l(f) \circ t$ and $l'(f') \circ t'$,

$$d(t \circ l(f), t' \circ l'(f')) = \min \begin{cases} d(t \circ f, t' \circ l'(f')) + \delta(l, \emptyset) \\ d(t \circ l(f), t' \circ f') + \delta(\emptyset, l') \\ d(t \circ f, t' \circ f') + \delta(l, l') \end{cases}$$



Rightmost decomposition

Substitution of l into l'



$\text{Distance}(l(f), l'(f'))$

+

$\text{Distance}(t, t')$



Deletion of l



$\text{Distance}(t \circ f, t' \circ l'(f')) + \text{del}(l)$

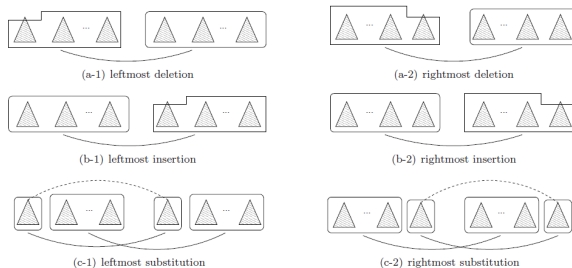
Insertion of l'



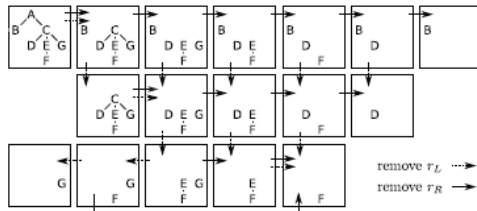
$\text{Distance}(t \circ l(f), t' \circ f') + \text{ins}(l')$

Relevant Sub-forests and Full Decomposition

Relevant sub-forests



The full decomposition

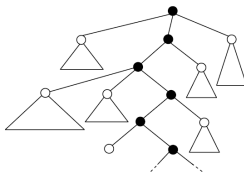


Root-leaf Path and Relevant Sub-trees

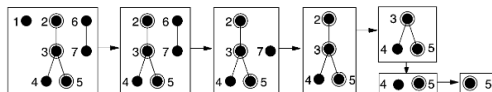
- Root-leaf path



- Relevant sub-trees

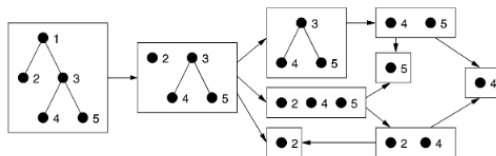


- Path decomposition

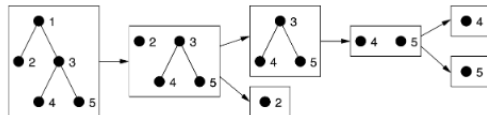


Path Decomposition

- Leftmost path decomposition



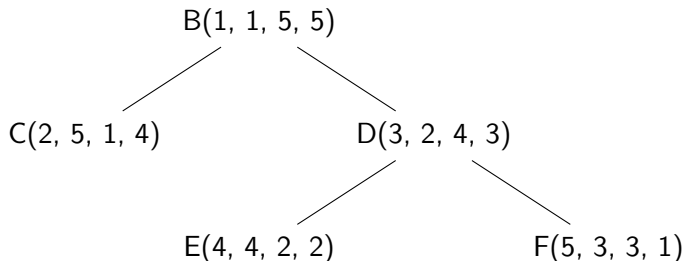
- Rightmost path decomposition



Tree Indexing

- $preOrder_left_to_right \rightarrow preL$
- $preOrder_right_to_left \rightarrow preR$
- $postOrder_left_to_right \rightarrow postL$
- $postOrder_right_to_left \rightarrow postR$

Consider a tree $T(preL, preR, postL, postR)$



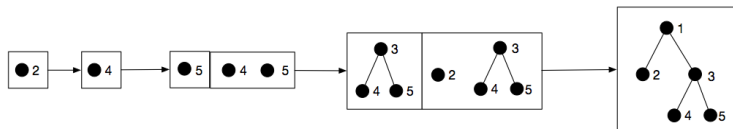
Note:

- $preL + postR = treeSize$
- $preR + postL = treeSize$

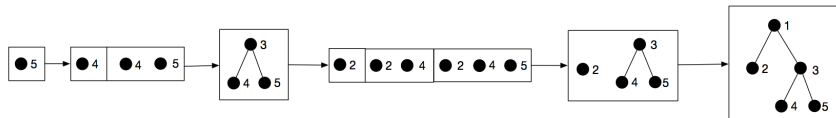
Bottom-up Enumeration

Two orders of the enumeration of the path decomposition:

- LR-postorder enumeration



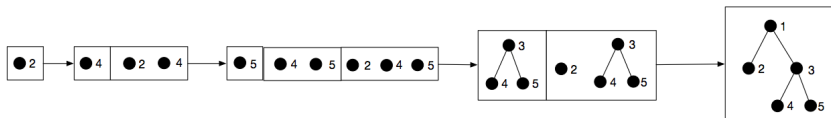
- RL-postorder enumeration



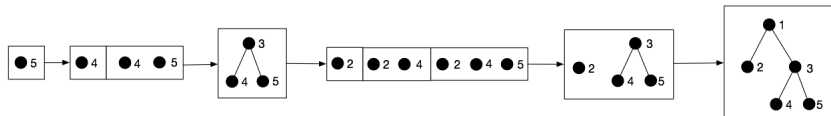
Bottom-up Enumeration

Two orders of the enumeration of the full decomposition:

- Prefix-suffix order



- Suffix-prefix order



A Simple Algorithm

Main Idea:

Enumerate each pairs of sub-forests in tree A and B in prefix-suffix post order or suffix-prefix post order.

Pseudocode:

Algorithm 1: Compute tree edit distance by enumerating all pairs in $O(m^2n^2)$ time.

inputs: (T_1, T_2) , with $|T_1| = m$ and $|T_2| = n$

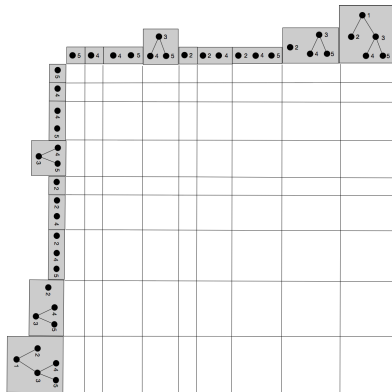
output: $d(T_1[i], T_2[j])$ for $1 \leq i \leq m$ and $1 \leq j \leq n$

```
1  $L_1 \leftarrow \text{POSTORDER}(T_1);$ 
2  $L_2 \leftarrow \text{POSTORDER}(T_2);$ 
3 for  $i = 1$  to  $|L_1|$  do
4   for  $j = 1$  to  $|L_2|$  do
5      $\mid$  compute  $d(L_1[i], L_2[j])$  as in Equation 2.3
6   end
7 end
8  $d(T_1[i], T_2[j]) \leftarrow d(L_1[|L_1|][|L_2|]);$ 
9 return  $d(T_1[i], T_2[j]);$ 
```

A Simple Algorithm

Implementation:

- Implemented in dynamic programming.
- Runs in $\mathcal{O}(m^2 n^2)$ time using $\mathcal{O}(m^2 n^2)$ space.

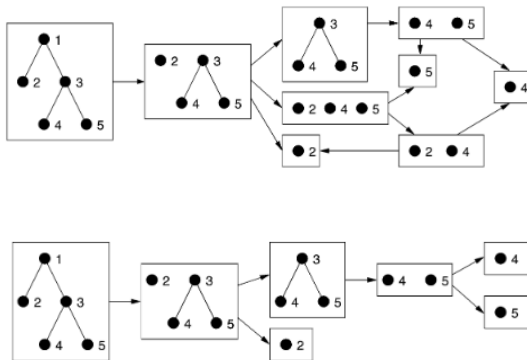


- The upper bound of the post-order enumeration method for the tree edit distance problem is $\mathcal{O}(n^4)$

Improved Algorithmic Path Strategies

Key Ideas:

- Different path decomposition creates different number of relevant sub-forests.
- Take advantage of the overlap among sub-forests that are contained in the same sub-tree, and the overlap of sub-trees as well.



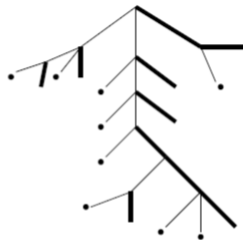
Zhang and Shasha's Algorithm

Key Ideas:

- Fixed direction in each recursive calls.
 - Recursive right decomposition - leftmost paths.
 - Recursive left decomposition - rightmost paths.



(a) leftmost paths



(b) rightmost paths

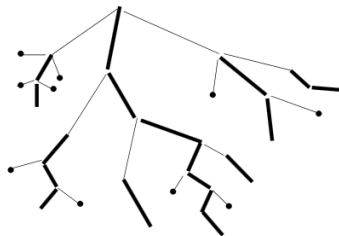
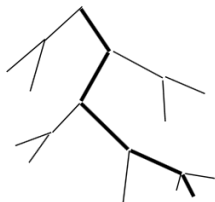
Implementation:

- Implemented via dynamic programming.
- Enumerates sub-trees pairs rooted at key roots.
 - Leftmost paths: either tree root or has a left sibling
 - Rightmost paths: either tree root or has a right sibling
- Uses temporary tables to store forest-forest distance and a permanent tables to store tree-tree distance.
- Runs in
 $\mathcal{O}(|A| * |B| * \min\{\text{depth}(A), \text{leaves}(A)\} * \min\{\text{depth}(B), \text{leaves}(B)\})$
time and $\mathcal{O}(|A| * |B|)$ space.

Klein's Algorithm

Key Ideas:

- Zhang and Shasha's algorithm depends on shapes of trees.
- Heavy paths generate the least number of sub-forests on one tree.



Implementation:

- Implemented via dynamic programming.
- applies heavy paths to the larger tree and fully decomposes the smaller tree.
- Uses temporary tables to store forest-forest distance and a permanent tables to store tree-tree distance.
- Runs in $\mathcal{O}(\log(|A|) * |B|^2)$ time using $\mathcal{O}(|A| * |B|)$ space.

Key Ideas:

- Klein's algorithm consider the shape of one tree with no consideration on the other tree.
- Heavy paths can be applied on both trees.
- Always applies heavy paths on the larger tree.

Implementation:

- Implemented via dynamic programming.
- Runs in $\mathcal{O}(|A|^2 |B| (1 + \log(\frac{|B|}{|A|})))$ time using $\mathcal{O}(|A| * |B|)$ space.

Conclusion

	Time	Space	Comments
Simple Algorithm	$\mathcal{O}(n^4)$	$\mathcal{O}(n^4)$	first algorithm
Zhang	$\mathcal{O}(n^2 \log^2(n))$	$\mathcal{O}(n^2)$	efficient for balanced trees
Klein	$\mathcal{O}(n^3 \log(n))$	$\mathcal{O}(n^2)$	no consider smaller trees
Demaine	$\mathcal{O}(n^3)$	$\mathcal{O}(n^2)$	worse case is frequent

Table: State-of-the-Art Algorithms in the Tree Edit Distance Problem

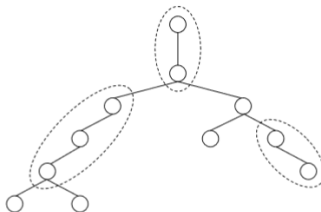
Key Idea:

- heavy paths is a kind of greedy algorithm which usually lead to the local optimum rather than global optimum.
- quantifies each possible paths and selects the path with the least number of sub-problems is selected
- implemented via dynamic programming.
- Preprocessing time is bounded by $\mathcal{O}(|A| |B|)$

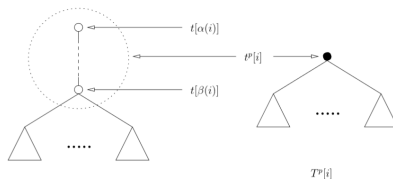
Vertical Reduction on Trees

Key Ideas:

- A large number of non-branching nodes can be compressed into one node.



- The compact representation of the original tree can aid in improving the running time for computing the tree edit distance.

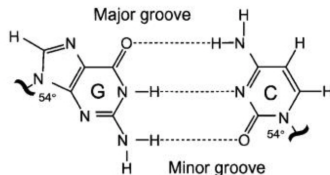
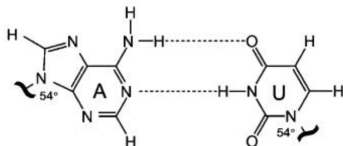


- RNA is an essential molecule in organisms.
 - Cellular organisms
 - Viruses
- RNA is assembled as a chain of nucleotides.

```

1  AAAGCAGGCC AGGCAACCGC UGCCUGCACC GCAAGGUGCA GGGGGAGGAA
51  AGUCCGGACU CCACAGGGCA GGGUGUUGGC UAACAGCCAU CCACGGCAAC
101 GUGCGGAAUA GGGCCACAGA GACGAGUCUU GCCGCCGGGU UGCCCCGGCG
151 GGAAGGGUGA AACGCGGUAA CCUCCACCUG GAGCAAUCCC AAAUAGGCAG
201 GCGAUGAAGC GGCCCGCUGA GUCUGCGGGU AGGGAGCUGG AGCCGGCUGG
251 UAACAGCCGG CCUAGAGGAA UGGUUGUCAC GCACCGUUUG CCGCAAGGCG
301 GGCGGGGCGC ACAGAAUCCG GCUUAUCGGC CUGCUUUGCU U
    
```

- RNA is unstable.



RNA Secondary Structure

- a list of base pairs satisfies the following constraints:
 - A base cannot participate in more than one base pair,
 - No two base pairs cross.

(1)	1	337	10
(2)	11	326	1
(3)	12	278	7
(4)	20	45	2
(5)	23	42	8
(6)	59	183	4
(7)	71	179	5
(8)	77	89	4
(9)	91	105	1
(10)	92	103	4
(11)	106	174	2
(12)	111	172	2
(13)	127	156	4
(14)	132	151	8
(15)	187	235	4
(16)	197	226	6
(17)	206	220	5
(18)	242	261	8
(19)	281	308	2
(20)	284	305	9

Tree Representation of the RNA Secondary Structure

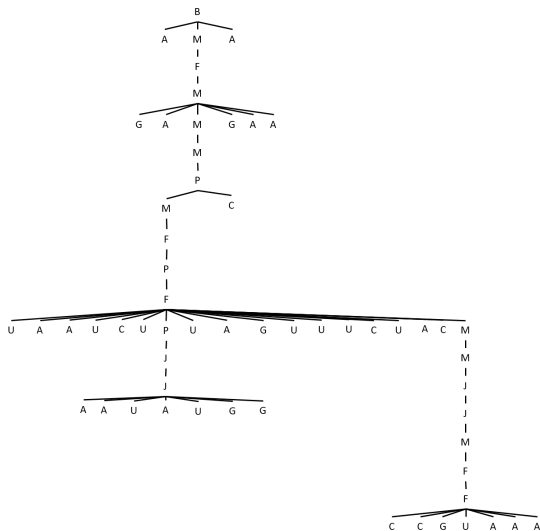
- Base pair labelling

Base Pairs	Label
$A = U$	F
$G \equiv C$	J
$U = A$	P
$C \equiv G$	M

Table: The Label of Base Pair

- Tree representation

Tree Representation of the RNA Secondary Structure



Experiment Result

RNAName: Alcaligenes-eutrophus-pb-b

RNASize: 343

```
((((((((-----((-----)))))--))-----((((-----((((-----)
AAAGCAGGCCAGGCAACCGUGCCUGCACCGCAAGGUGCAGGGGGAGGAAAGUCCGGACUCCACAGGGCAGGGUGUUGGCUAACAG
)))-----))--((-----((((-----)))))--))-----((((-----)
CCAUCCACGGCAACGUGCGGAAUAGGGCCACAGAGACGAGUCUUGCCGCCGGGUUCGCCGGCGGGAAGGGUGAAA-----
)))))--))--((-----))--))-----))--))-----))--((-----))--))-----))--((-----))--))
GAGUCUGCGGGUAGGGAGCUGGAGCCGGCUGGUAACAGCCGGCCUAGAGGAAUGGUUGUCACGCACCGUUUGCCGCAAGGCGGGC
)-))-----))--))--))--))-----
GGGCGCACAGAAUCCGGCUUAUCGGCCUGCUUUGCUU
```

RNAName: Streptomyces-bikiniensis-gpb-h

RNASize: 400

```
((((((((-----((-----)))))--))-----((((-----((((-----)
CGAGCCGGGCGGGCGGCCGCGUGGGGGUCUUCGACCUCGCCGAGGAACGUCCGGGCUCCACAGAGCAGGGUGGUGGCUAACGGCCACCCGG
(-----))--((-----((((-----)))))--))-----((((-----)
GGUGACCCGCGGGACAGUGCCACAGAAAACAGACCGCCGGGACCUCGGUCCUCGGUAAGGGUGAAACGGUGGUGUAAGAGACCACCAGC
)))))--))--((-----))--))-----))--))-----))--((-----))--))
CGAGUCCGCGGGUAGACCGCACGAGGCCGGCGGCAACGCCGCCUAGAUGGAUGGCCGUCGCCCGACGACCGCGAGGUCCCGGGG
-----))--))--))--))-----
ACAGAACCCGGCGUACAGCCCGACUCGUCUG
```

Time Complexity(Naive Algorithm)

Prerequisite of the Analysis

- Lemma 1 implies dynamic programming.
- Each new subproblem can be computed in constant time.
- Time complexity is bounded by
 $\#subproblems\ of\ F_1 * \#subproblems\ of\ F_2.$

Count the subproblems of F

- $\#subproblem\ of\ F = \#(i,j) - deleted\ subforest, 0 \leq i + j \leq |F|$
- $\#(i,j) - deleted\ subforest = \sum_{k=0}^{|F|} k = \mathcal{O}(|F|^2)$, since for each i there are $|F| - i$ choices for j
- numerous redundant subforests.

Conclusion

Time complexity is bounded by $\mathcal{O}(|F_1|^2 |F_2|^2)$.

Key Idea(Zhang and Shasha's Algorithm)

LEMMA 2

For tree-tree edit distance

$$\delta(l(f), l'(f')) = \min \begin{cases} \delta(f, l'(f')) + \gamma(l \rightarrow \lambda) \\ \delta(l(f), f') + \gamma(\lambda \rightarrow l') \\ \delta(f, f') + \gamma(l \rightarrow l') \end{cases} \quad (1)$$

For forest-forest edit distance

$$\delta(l(f) \circ t, l'(f') \circ t') = \min \begin{cases} \delta(f \circ t, l'(f') \circ t') + \gamma(l \rightarrow \lambda) \\ \delta(l(f) \circ t, f' \circ t') + \gamma(\lambda \rightarrow l') \\ \delta(l(f), l'(f')) + \delta(t, t') \end{cases} \quad (2)$$

Key Idea(Zhang and Shasha's Algorithm)

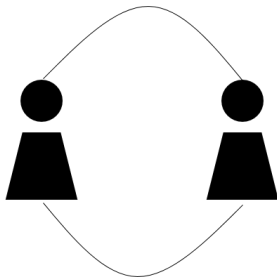
The lemma 2 is correct.

Proof.

For tree-tree distance, formula

$$\delta(l(f), l'(f')) \leq \delta(f, f') + \gamma(l, l')$$

represents a particular mapping of $\delta(l(f), l'(f'))$



Key Idea(Zhang and Shasha's Algorithm)

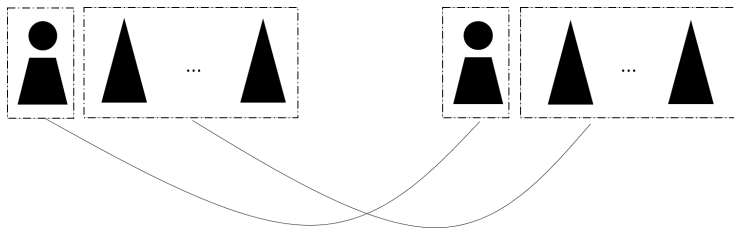
The lemma 2 is correct.

Proof.Continue

For forest-forest distance, formula

$$\delta(l(f) \circ t, l'(f') \circ t') \leq \delta(l(f), l'(f')) + \delta(t, t')$$

represents a particular mapping of $\delta(l(f) \circ t, l'(f') \circ t')$



Key Idea(Zhang and Shasha's Algorithm)

The lemma 2 is correct.

Another Proof.

Naive Algorithm

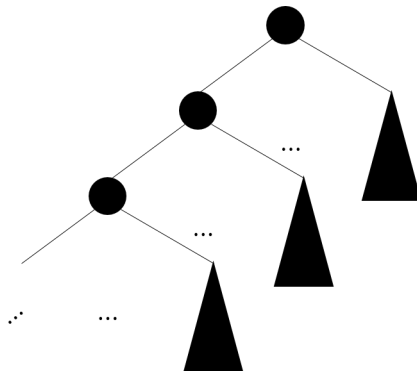
$$\delta(l(f) \circ t, l'(f') \circ t') = \min \begin{cases} \delta(f \circ t, l'(f') \circ t') + \gamma(l \rightarrow \lambda) \\ \delta(l(f) \circ t, f' \circ t') + \gamma(\lambda \rightarrow l') \\ \delta(f, f') + \delta(t, t') + \gamma(l \rightarrow l') \end{cases}$$

Zhang-Shasha's Algorithm

$$\delta(l(f) \circ t, l'(f') \circ t') = \min \begin{cases} \delta(f \circ t, l'(f') \circ t') + \gamma(l \rightarrow \lambda) \\ \delta(l(f) \circ t, f' \circ t') + \gamma(\lambda \rightarrow l') \\ \delta(l(f), l'(f')) + \delta(t, t') \end{cases}$$

Key Idea(Zhang and Shasha's Algorithm)

- dynamic programming style algorithm
- bottom-up procedure
- compute tree-tree distance between subtrees rooted at keyroots of both tree to avoid redundant computation.



Key Concept(Zhang and Shasha's Algorithm)

We define keyroot of T as

$$\text{keyroots}(T) = \{k \mid k \text{ is either root}(T) \text{ or } k \text{ has at least one siblings}\}.$$

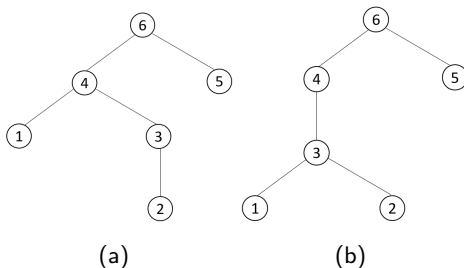


Figure: compute the edit distance between tree(a) and tree(b)

$$\text{keyroot}(T_1) = \{3, 5, 6\}$$

$$\text{keyroot}(T_2) = \{2, 5, 6\}$$

Time Complexity(Zhang and Shasha's Algorithm)

Prerequisite of the Analysis

- The time complexity is bounded by

$$\sum_{i=1}^{|keyroots(T_1)|} |T_1(keyroots(T_1)[i])| * \sum_{i=1}^{|keyroots(T_2)|} |T_2(keyroots(T_2)[i])|$$

Time Complexity(Zhang and Shasha's Algorithm)

Prerequisite of the Analysis

- Define for each node i ,

$$cdepth(i) = |anc(i) \cap keyroots(T)|$$

Define for a tree T ,

$$cdepth(T) = \max cdepth(i)$$

- **Lemma 3**

$$|keyroot(T)| \leq |leaves(T)|$$

, since each leaf is the leftmost descendant of at most one member of $keyroot(T)$.

- From definition and Lemma 3 $cdepth(i) \leq \min(depth(T), leaves(T))$ for $1 \leq i \leq |T|$. Hence, $cdepth(T) \leq \min(depth(T), leaves(T))$.

Time Complexity(Zhang and Shasha's Algorithm)

Prerequisite of the Analysis

• Lemma 4

$$\sum_{i=1}^{|keyroots(T)|} |T(keyroots(T)[i])| = \sum_{j=1}^{|T|} cdepth(j)$$

, each node j is counted $cdepth(j)$ times.

Time Complexity(Zhang and Shasha's Algorithm)

Conclusion

- Time Complexity

$$\begin{aligned} & i=|keyroots(T_1)| \quad j=|keyroots(T_2)| \\ & \sum_{i=1} \quad \sum_{j=1} \quad T_1(keyroots(T_1)[i]) * T_2(keyroots(T_2)[j]) \\ & = \sum_{i=1}^{|T_1|} cdepth(i) * \sum_{j=1}^{|T_2|} cdepth(j) \\ & \leq |T_1| * cdepth(T_1) * |T_2| * cdepth(T_2) \end{aligned}$$

- Hence we can conclude that the time complexity is

$$\mathcal{O}(|T_1| * |T_2| * \min(depth(T_1), leaves(T_1)) * \min(depth(T_2), leaves(T_2)))$$

Definition(Klein's Algorithm)

Leftmost decomposition

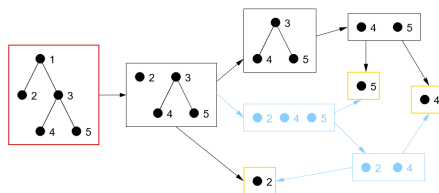
$$\delta(t \circ l(f), t' \circ l'(f')) = \min \begin{cases} \gamma(l \rightarrow \lambda) + \delta(t \circ f, t' \circ l'(f')) \\ \gamma(\lambda \rightarrow l') + \delta(t \circ l(f), t' \circ f') \\ \delta(l(f), l'(f')) + \delta(t, t') \end{cases}$$

Rightmost decomposition

$$\delta(l(f) \circ t, l'(f') \circ t') = \min \begin{cases} \gamma(l \rightarrow \lambda) + \delta(f \circ t, l'(f') \circ t') \\ \gamma(\lambda \rightarrow l') + \delta(l(f) \circ t, f' \circ t') \\ \delta(l(f), l'(f')) + \delta(t, t') \end{cases}$$

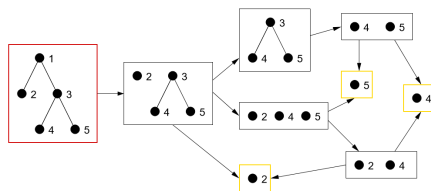
Observation(Klein's Algorithm)

Left Decomposition



gives 7 sub-forests.

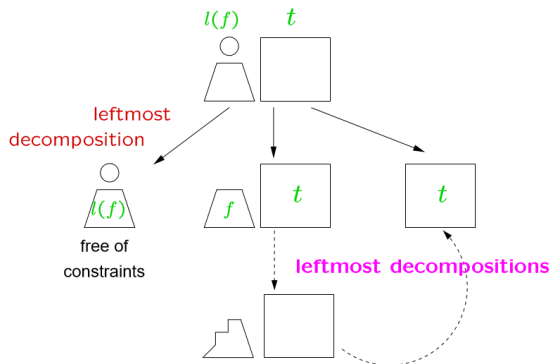
Right Decomposition



gives 9 sub-forests.

Key Idea(Klein's Algorithm)

- Optimize the subproblems of a forest($l(f) \circ t$) by eliminating nodes of f in $f \circ t$, so that $f \circ t$ and t share relevant forests as most as possible.
- However, this method also loses the ability to control the second forest($l'(f') \circ t'$)



Key Concept(Klein's Algorithm)

- **heavy child.** $heavy(i)$ denote the child of i having greatest weight.
- **heavy path.** The sequence of nodes $i, heavy(i), heavy(heavy(i)), \dots$ defines a descending path which is called the heavy path.

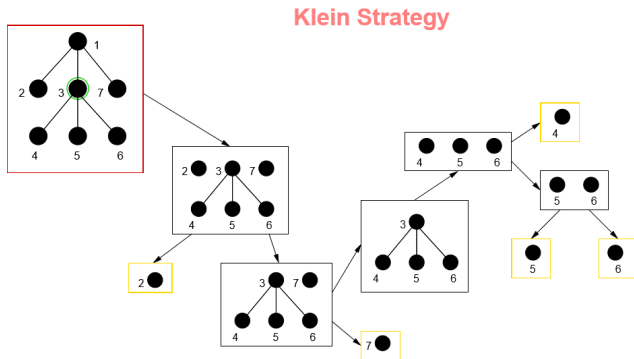


Algorithm(Klein's Algorithm)

uses the heavy path as a guide for the decomposition of the forest.

Let $I(f) \circ t$ be the forest to be decomposed. compute the heavy path P for $I(f) \circ t$

- 1) if I belongs to P , apply $t \circ I(f)$. Otherwise, apply $I(f) \circ t$.
- 2) apply this scheme recursively to all sub-forests of $I(f) \circ t$.



Time Complexity(Klein's Algorithm)

Prerequisite of the Analysis

- Define one of node i 's children with maximal number of descendants and marks it as *heavy*, and we mark all the other children of i as *light*
- The time complexity is bounded by

$$\sum_{i \in \text{light}(T_1)} |T_1(i)| * |T_2|^2$$

Time Complexity(Klein's Algorithm)

Prerequisite of the Analysis

- Define $ldepth(i)$ to be the number of light nodes that are **proper ancestors** of i in F .
- **Theorem 1**(Sleator and Tarjan, 1983)
Let v be any vertex of tree F . There is at most one heavy edge entering v , and there are at most $\log_2 |F|$ light edges on the tree path from v to $root(v)$.
- According to theorem 1, for each node i in tree F , we have

$$ldepth(i) \leq \log_2 |F|$$

Time Complexity(Klein's Algorithm)

Conclusion

- Subforests of T_1

$$\begin{aligned} & \sum_{i \in \text{light}(T_1)} |T_1(i)| \\ & \leq \sum_{i \in T_1} (1 + \text{ldepth}(i)) \\ & \leq \sum_{i \in T_1} (\log_2 |T_1| + \mathcal{O}(1)) \\ & = \mathcal{O}(|T_1| \log_2 |T_1|) \end{aligned}$$

- Hence the time complexity is $\mathcal{O}(|T_1| \log_2 |T_1| * |T_2|^2)$

Lemma 5

- if F or F' is an empty forest:

$R(\epsilon, l(f) \circ t) = (\epsilon, l(f) \circ t) \cup R(\epsilon, f \circ t)$ whenever $\phi(\epsilon, l(f) \circ t) = \text{left}$,
 $R(\epsilon, t \circ l(f)) = (\epsilon, t \circ l(f)) \cup R(\epsilon, t \circ f)$ whenever $\phi(\epsilon, t \circ l(f)) = \text{right}$,

$R(l(f) \circ t, \epsilon) = (l(f) \circ t, \epsilon) \cup R(f \circ t, \epsilon)$ whenever $\phi(l(f) \circ t, \epsilon) = \text{left}$,
 $R(t \circ l(f), \epsilon) = (t \circ l(f), \epsilon) \cup R(t \circ f, \epsilon)$ whenever $\phi(t \circ l(f), \epsilon) = \text{right}$,

- if $(F, F') = (l(f) \circ t, l'(f') \circ t')$ and $\phi(F, F') = \text{left}$, then
 $R(F, F') = (F, F') \cup R(f \circ t, F') \cup R(F, f' \circ t') \cup R(l(f), l'(f')) \cup R(t, t')$
- if $(F, F') = (t \circ l(f), t' \circ l'(f'))$ and $\phi(F, F') = \text{right}$, then
 $R(F, F') = (F, F') \cup R(t \circ f, F') \cup R(F, t \circ f) \cup R(l(f), l'(f')) \cup R(t, t')$

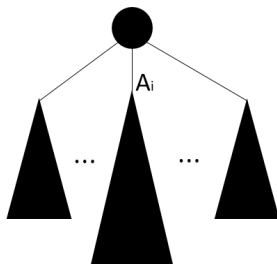
Lemma 6

- $R(I(f)) = \{I(f)\} \cup R(f)$, no matter what the direction is.
That is to say, $\#ref(I(f)) \geq |I(f)|$
- $R(I(f) \circ t) = \{(I(f) \circ t)\} \cup R(f \circ t) \cup R(I(f)) \cup R(t)$, if the direction is left.
That is to say, $\#ref(I(f) \circ t) \geq \#ref(I(f)) + \#ref(t)$
- $R(t \circ I(f)) = \{(t \circ I(f))\} \cup R(t \circ f) \cup R(I(f)) \cup R(t)$, if the direction is right.

Lower Bound of Relevant Forests

Lemma 7

Given a tree $A = I(A_1 \circ A_2 \circ \dots \circ A_n)$, for any strategy we have,
 $\#rel(A) \geq |A| - |A_i| + \#rel(A_1) + \dots + \#rel(A_n)$
, where $i \in [1 \dots n]$ is such that size of A_i is maximal.



Lemma 7

Given a tree $A = I(A_1 \circ A_2 \circ \dots \circ A_n)$, for any strategy we have,
 $\#rel(A) \geq |A| - |A_i| + \#rel(A_1) + \dots + \#rel(A_n)$
, where $i \in [1 \dots n]$ is such that size of A_i is maximal.

Proof.

From lemma 6, $R(I(f)) = I(f) \cup R(f)$. That is to say,
 $\#rel(I(f)) = 1 + \#rel(f)$.

Then let $F = A_1 \circ \dots \circ A_n$. To prove lemma 7, we first prove that
 $\#ref(F) \geq |F| - |A_i| + \#ref(A_1) + \dots + \#ref(A_n)$

Lower Bound of Relevant Forests

Lemma 7

Given a tree $A = l(A_1 \circ A_2 \circ \dots \circ A_n)$, for any strategy we have,
 $\#rel(A) \geq |A| - |A_i| + \#rel(A_1) + \dots + \#rel(A_n)$
, where $i \in [1 \dots n]$ is such that size of A_i is maximal.

Proof.

For case $n = 1$, $\#ref(A_1) \geq |A_1| - |A_1| + \#ref(A_1)$

For case $n > 1$, let $F = l(g) \circ A_2 \dots A_n$, **t** be $A_2 \circ \dots \circ A_n$ then by lemma 5,

$$R(F) = \{F\} \cup R(g \circ t) \cup R(A_1) \cup R(t)$$

From lemma 6, we know,

$$R(g \circ t) \geq |g| + |t| \geq \min\{|g|, |t|\}$$

By applying the inequality, we have

$$\#ref(F) \geq 1 + \#ref(A_1) + \#ref(t) + \min\{|g|, |t|\}$$

Lower Bound of Relevant Forests

Lemma 7

Given a tree $A = I(A_1 \circ A_2 \circ \dots \circ A_n)$, for any strategy we have,
 $\#rel(A) \geq |A| - |A_i| + \#rel(A_1) + \dots + \#rel(A_n)$
, where $i \in [1 \dots n]$ is such that size of A_i is maximal.

Proof.

Applying induction hypothesis for $t = A_2 \circ \dots \circ A_n$
 $\#ref(F) \geq 1 + \#ref(A_1) + \dots + \#ref(t) + |t| - |A_j| + \min\{|g|, |t|\}$

Then we need to prove that, $1 + |t| - |A_j| + \min\{|g|, |t|\} \geq |F| - |A_i|$,
where A_i is the size of the largest sub-tree in forest $A_1 \circ \dots \circ A_n$, while A_j
is the size of the largest sub-tree in forest $A_2 \circ \dots \circ A_n$

Lemma 7

Given a tree $A = I(A_1 \circ A_2 \circ \dots \circ A_n)$, for any strategy we have,
 $\#rel(A) \geq |A| - |A_i| + \#rel(A_1) + \dots + \#rel(A_n)$
, where $i \in [1 \dots n]$ is such that size of A_i is maximal.

Proof.

In the case when $|g| \leq |t|$, then we have $1 + |t| + |g| = |F|$. Since $|A_j| \leq |A_i|$, it follows that $1 + |t| - |A_j| + \min\{|g|, |t|\} \geq |F| - |A_i|$

In the case when $|g| > |t|$, this implies that A_1 is the largest subtree of A .
So $i = 1$ and $|F| - |A_i| = |t|$, that proves
 $1 + |t| - |A_j| + \min\{|g|, |t|\} \geq |F| - |A_i|$

Lower Bound of Relevant Forests

Lemma 7

Given a tree $A = l(A_1 \circ A_2 \circ \dots \circ A_n)$, for any strategy we have,
 $\#rel(A) \geq |A| - |A_i| + \#rel(A_1) + \dots + \#rel(A_n)$
, where $i \in [1 \dots n]$ is such that size of A_i is maximal.

Proof.

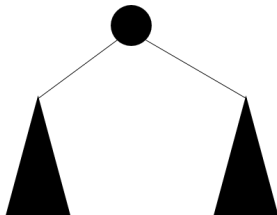
To sum up, the number of relevant sub-forest of a tree A is

$$\begin{aligned}\#ref(A) &= 1 + \#rel(F) \\ &\geq 1 + |F| - |A_i| + \#rel(A_1) + \dots + \#rel(A_n) \\ &= |A| - |A_i| + \#rel(A_1) + \dots + \#rel(A_n)\end{aligned}$$

Lower Bound of Relevant Forests

Lemma 8 $\#rel(A)$ has a lower bound in $\mathcal{O}(n \log(n))$

Proof. Let T_n be a complete balanced binary tree of size n , we prove that $\#rel(T_n) \geq \frac{(n+1)\log_2(n+1)}{2}$



For case $n = 1$, $rel(T_n) = 1$, consistent with lemma 8.

Lower Bound of Relevant Forests

Lemma 8 $\#rel(A)$ has a lower bound in $\mathcal{O}(n \log(n))$

Proof. Let T_n be a complete balanced binary tree of size n , we prove that $\#ref(T_n) \geq \frac{(n+1)\log_2(n+1)}{2}$

For case $n > 1$, by lemma 7, we have

$$\#rel(T_n) \geq n - m + 2 * \#ref(T_m)$$

, where $m = (n - 1)/2$. T_n is of the form $l(A \circ B)$, where A and B are two complete balanced tree of size m . By induction hypothesis for T_m , it converts to

$$\begin{aligned} \#ref(T_n) &\geq n - m + (m + 1)\log_2(m + 1) \\ &= \frac{(n + 1)(\log_2(m + 1) + 1)}{2} \\ &= \frac{(n + 1)(\log_2(2m + 1) + 1)}{2} \\ &= \frac{(n + 1)(\log_2(n + 1))}{2} \end{aligned}$$

Lemma 9

Let A and B be two trees of size n . For any decomposition strategy, $\#rel(A, B)$ has a lower bound in $\mathcal{O}(n^2 \log^2(n))$

Cover Strategies

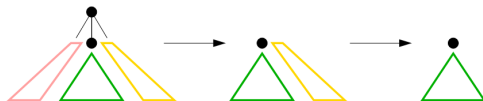
Cover

Let F be a tree, for each node i in F

- if $\deg(i) = 0$ or $\deg(i) = 1$, then $r(i) \in \{right, left\}$.
- if $\deg(i) > 1$, then $r(i)$ is a favourite child of i .

Cover Strategies

- if $\deg(i) = 0$ or $\deg(i) = 1$, then $\phi(A(i), G) = r(i)$, for each forest G of B .
- Otherwise, let $A(i)$ is of the form $I(A_1 \circ \dots \circ A_p \circ \dots \circ A_n)$
 - $\phi(T \circ A_p \dots \circ A_n, G) = left$, for each forest T of $A_1 \circ \dots \circ A_{p-1}$
 - $\phi(A_p \circ T, G) = right$, for each forest T of $A_{p+1} \circ \dots \circ A_n$



leftmost decompositions

rightmost decompositions

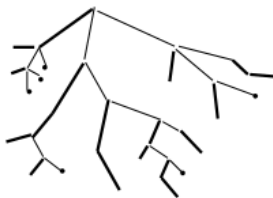
Cover Strategies

Zhang-Shasha Algorithms Cover Strategies

- for any node of degree 0 or 1, the direction is $\text{right}(t \circ l(f))$.
- for any other node, the favorite child is the leftmost child. That is to say, the direction is $\text{right}(t \circ l(f))$.

Klein's Algorithm Cover Strategies

- for any node of degree 0 or 1, the direction is $\text{left}(l(f) \circ t)$.
- for any other node, the favorite child is the root of the heaviest subtree.



(a)



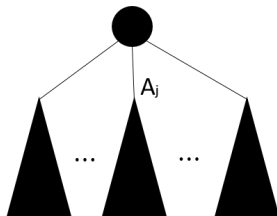
(b)

Exact Number of Relevant Forests

Lemma 10

Let $A = I(A_1 \circ \cdots \circ A_n)$ be a cover tree such that $n = 1$ or the root of A_j is the favourite child.

$$\#rel(A) = |A| - |A_j| + \#rel(A_1) + \cdots + \#rel(A_n)$$



Upper Bound of Relevant Forests

Lemma 11

Let $A = I(A_1 \circ \cdots \circ A_n)$ be a cover tree

$$\#rel(A) \leq \frac{n(n+3)}{2} - \sum_{i \in A} |A(i)|$$

In average

$$\frac{1}{2}n^2 + \frac{\sqrt{\pi}}{2}n^{\frac{3}{2}} + \mathcal{O}(n)$$

Happens of a Cover Strategy on a tree to Another Tree

Lemma 12

Given a pair of tree (A, B) provided with a cover for A , all relevant forests of A fall within three categories:

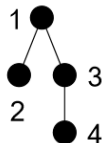
- α those that are compared with all rightmost forests of B ,
- β those that are compared with all leftmost forests of B ,
- γ those that are compared with all forests of B .

- $\#left(B)$: number of leftmost subforests of B
- $\#right(B)$: number of rightmost subforests of B
- $\#special(B)$: number of all subforests of B

Note: $\#left(B)$, $\#right(B)$, $\#special(B)$ are known, depended by the structure of B

Happens of a Cover Strategy on a tree to Another Tree

- $\#left(B)$: number of leftmost subforests of B
- $\#right(B)$: number of rightmost subforests of B
- $\#special(B)$: number of all subforests of B



$$\#right(B) = 4 - 2 + 1 + 2 = 5$$

$$\#left(B) = 4 - 1 + 1 + 2 = 6$$

$$\#special(B) = 4 * 7/2 - 4 - 1 - 2 - 1 = 6$$

Status of Nodes in a Cover Tree

- **Free** : nodes that do not receive anything from the parent
- **Left** : nodes that inherit leftmost forests of B.

$$Left(A(i)) = (A(i), \text{all leftmost forests of } B)$$

- **Right** : nodes that inherit rightmost forests of B.

$$Right(A(i)) = (A(i), \text{all rightmost forests of } B)$$

- **All** : nodes that inherit all subforests of B from the parent

$$All(A(i)) = (A(i), \text{all subforests of } B)$$

Note1: The status of a node depends of the direction and of the heritage from parent.

Note2: By definition, $\#rel(A, B)$ equals $Free(A)$.

LEMMA 13

Let (A, B) be a pair of trees. A being a cover tree.

- 1) If A is reduced to a single node whose direction is right.
 - $\text{Free}(A) = \text{Left}(A) = \#left(B)$
 - $\text{All}(A) = \text{Right}(A) = \#special(B)$
- 2) If A is reduced to a single node whose direction is left.
 - $\text{Free}(A) = \text{Right}(A) = \#right(B)$
 - $\text{All}(A) = \text{Left}(A) = \#special(B)$

Number of Relevant Forests of A Pair of Tree

LEMMA 13

Let (A, B) be a pair of trees. A being a cover tree.

- 3) If $A = I(A')$ and the direction of I is right
- $\text{Free}(A) = \text{Left}(A) = \# \text{left}(B) + \text{Right}(A')$
?? $\text{Free}(A) = \text{Left}(A) = \# \text{left}(B) + \text{Left}(A')$
 - $\text{All}(A) = \text{Right}(A) = \# \text{special}(B) + \text{All}(A')$
- 4) If $A = I(A')$ and the direction of I is left
- $\text{Free}(A) = \text{Right}(A) = \# \text{right}(B) + \text{Left}(A')$
?? $\text{Free}(A) = \text{Right}(A) = \# \text{left}(B) + \text{Right}(A')$
 - $\text{All}(A) = \text{Left}(A) = \# \text{special}(B) + \text{All}(A')$



Number of Relevant Forests of A Pair of Tree

LEMMA 13

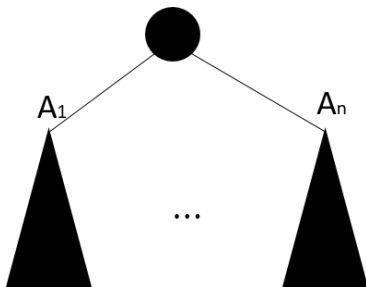
Let (A, B) be a pair of trees. A being a cover tree.

5) If $A = I(A_1 \circ \dots \circ A_n)$ and the favourite child is A_1

- $\text{Free}(A) = \text{Left}(A) = \sum_{i>1} \text{Free}(A_i) + \text{Left}(A_1) + \# \text{left}(B)(|A| - |A_1|)$
- $\text{All}(A) = \text{Right}(A) = \sum_{i>1} \text{Free}(A_i) + \text{All}(A_1) + \# \text{special}(B)(|A| - |A_1|)$

6) If $A = I(A_1 \circ \dots \circ A_n)$ and the favourite child is A_n

- $\text{Free}(A) = \text{Left}(A) = \sum_{i<n} \text{Free}(A_i) + \text{Left}(A_n) + \# \text{left}(B)(|A| - |A_n|)$
- $\text{All}(A) = \text{Right}(A) = \sum_{i<n} \text{Free}(A_i) + \text{All}(A_n) + \# \text{special}(B)(|A| - |A_n|)$



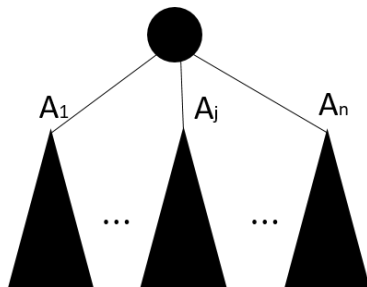
Number of Relevant Forests of A Pair of Tree

LEMMA 13

Let (A, B) be a pair of trees. A being a cover tree.

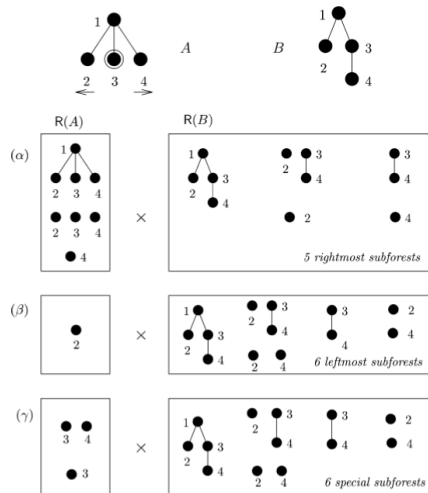
7) If $A = I(A_1 \circ \dots \circ A_n)$ and the favourite child is A_1

- $\text{Free}(A) = \text{Right}(A) = \sum_{i \neq j} \text{Free}(A_i) + \text{All}(A_j) + \# \text{right}(B)(1 + |A_1 \circ \dots \circ A_{j-1}|) + \# \text{special}(B) |A_{j+1} \circ \dots \circ A_n|$
- $\text{All}(A) = \text{Left}(A) = \sum_{i \neq j} \text{Free}(A_i) + \text{All}(A_j) + \# \text{special}(B)(|A| - |A_j|)$



Number of Relevant Forests of A Pair of Tree

Example of LEMMA 13



gives 33 pairs of subforest.

Example of LEMMA 13

$$\begin{aligned}Free(A) &= Free(2) + Free(4) + All(3) + \#special(B) * 1 + \#right(B) * 2 \\&= \#left(B) + \#right(B) + \#special(B) + \#special(B) \\&\quad + \#right(B) * 2 \\&= 33\end{aligned}$$

Number of Rightmost, Leftmost and Subforest.

LEMMA 14

Let A be a tree,

$$\#right(A) = \sum(|A(i)|, i \in A) - \sum(|A(j)|, j \text{ is a rightmost child})$$

$$\#left(A) = \sum(|A(i)|, i \in A) - \sum(|A(j)|, j \text{ is a leftmost child})$$

Proof.

From lemma 10, let A be $I(A_1 \circ \dots \circ A_n)$, then

$$\#right(A) = \sum_i \#right(A_i) + |A| - |A_n|$$

Number of Rightmost, Leftmost and Subforest.

LEMMA 15

Let F be a forest of size n ,

$$\#special(F) = \frac{n(n+3)}{2} - \sum_{i \in F} |F(i)|$$

Proof

For case $n = 0$, then $\#special(F) = 0$, which is consistent with lemma 15.

For case $n > 1$, let $F = l(g) \circ t$. There are two kinds of special sub-forest of F to be considered:

- those containing the node l : there are $|t| + 1$ such sub-forest
- those not containing the node l : there are $\#special(g \circ t)$ such sub-forest

Number of Rightmost, Leftmost and Subforest.

LEMMA 15

Let F be a forest of size n ,

$$\#special(F) = \frac{n(n+3)}{2} - \sum_{i \in F} |F(i)|$$

Proof

This allows that

$$\#special(F) = |t| + 1 + \#special(g \circ t)$$

On one hand, $|t| + 1 = n - |l(g)| + 1$

On the other hand, by applying the hypothesis on $g \circ t$, whose size is $n-1$

$$\#special(g \circ t) = \frac{(n-1)(n+2)}{2} - \sum_{i \in g \circ t} |g \circ t(i)|$$

Since $g \circ t$ is a sub-forest of F , this implies

$$\#special(g \circ t) = \frac{(n-1)(n+2)}{2} - \sum_{i \in g \circ t} |F(i)|$$

Number of Rightmost, Leftmost and Subforest.

LEMMA 15

Let F be a forest of size n ,

$$\#special(F) = \frac{n(n+3)}{2} - \sum_{i \in F} |F(i)|$$

Proof

It allows that

$$\begin{aligned}\#special(F) &= n - |l(g)| + 1 + \frac{(n-1)(n+2)}{2} - \sum_{i \in g \circ t} |F(i)| \\ &= n + 1 + \frac{(n-1)(n+2)}{2} - \sum_{i \in F} |F(i)| \\ &= \frac{n(n+3)}{2} - \sum_{i \in F} |F(i)|\end{aligned}$$

This concludes the proof.

LEMMA 16

For Zhang and Shasha's algorithm, $\#rel(A, B) = \#left(A) * \#left(B)$.

Proof

Applying lemma 13, case 1, 3, 5 and lemma 14. This concludes the proof.

Key Idea(A New Optimal Cover Strategy)

- An optimal cover strategy is a strategy that minimizes the total number of relevant forest.
- Define four dynamic programming tables **Right**, **Left**, **Free** and **All** indexed by nodes A.
- The favourite child is chosen to be the child that minimizes the number of relevant forests.

Algorithm(A New Optimal Cover Strategy)

Let $A = I(A_1 \circ \dots \circ A_n)$ then

$$\begin{aligned} \text{Free}(A) = & \sum_{i \geq 1} \text{Free}(A_i) \\ & + \min \begin{cases} \text{Left}(A_1) - \text{Free}(A_1) + \# \text{left}(B) * (|A| - |A_1|), \\ \text{All}(A_j) - \text{Free}(A_j) + \# \text{special}(B) |A_{j+1} \circ \dots \circ A_n| \\ \quad + \# \text{right}(B)(1 + |A_1 \circ \dots \circ A_{j-1}|), 1 < j < n \\ \text{Right}(A_n) - \text{Free}(A_n) + \# \text{right}(B)(|A| - |A_n|). \end{cases} \end{aligned}$$

Note 1 The favourite child is selected to be the root of the subtree A_j so that A_j gives the minimal value.

Note 2 The optimal cover is then built up by tracing back from $\text{Free}(A)$.

Time Complexity(A New Optimal Cover Strategy)

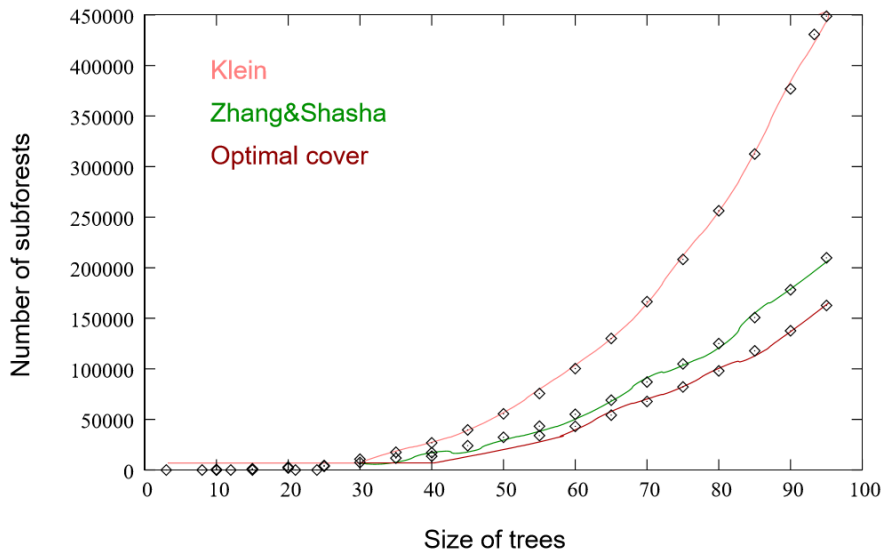
Prerequisite of the Analysis

- Two main steps in preprocessing
 - (1) the computation of $\#right(B)$, $\#left(B)$ and $\#special(B)$.
 - (2) the computation of array of **Right**, **Left**, **Free** and **All** indexed by nodes A.
- Time complexity is bounded by $\mathcal{O}((1) + (2))$
- (1) can be computed by applying lemma 10 and 11 in $\mathcal{O}(|B|)$
- Each array in (2) can be computed in $\mathcal{O}(\sum_{i \in A} deg(i))$, that is in $\mathcal{O}(|A|)$

Conclusion

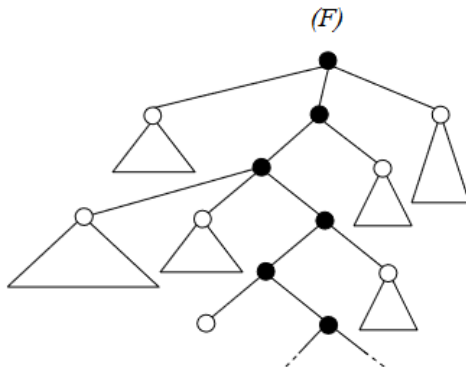
The time complexity of preprocessing is $\mathcal{O}(|A| + |B|)$

Performance(A New Optimal Cover Strategy)



Key Concept(An Optimal Decomposition Algorithm)

- set **TopLight(F)** is the set of *light nodes* with *ldepth* 1 in F.



Key Idea(An Optimal Decomposition Algorithm)

- For all $v \in \text{TopLight}(F)$, Klein's strategy solves $\delta(F(v), G)$ by determining the direction according to $F(v)$, even if $|F(v)| < |G|$
- In such case, making decision according to the larger forest can do better.

Algorithm(An Optimal Decomposition Algorithm)

We compute $\delta(F, G)$ recursively as follows:

- (1) if $|F| < |G|$, compute $\delta(G, F)$ instead.
- (2) Recursively compute $\delta(F(v), G)$ for all $v \in TopLight(F)$.
- (3) Compute $\delta(F, G)$ using the following decomposition strategy:
 - $S(F', G') = \text{left}$, if F' is a tree, or the root of leftmost tree is not the heavy child of its parent.
 - Otherwise, $S(F', G') = \text{right}$.

Time Complexity(An Optimal Decomposition Algorithm)

Prerequisite of the Analysis

- $\sum_{v \in TopLight(F)} |F(v)| \leq |F|$.
- $|F(v)| < \frac{|F|}{2}$ for each $v \in TopLight(F)$. Otherwise v would be a heavy node.
- This algorithm consists of two main steps(step 2 and 3). Then algorithm is bounded by $\mathcal{O}(\mathcal{O}(Step2) + \mathcal{O}(Step3))$
- In step 3, a node v in the heavy path of F cannot be matched or deleted until the remaining subforest is a tree. That is to say only $|F|$ subforest need to be considered. This gives $\mathcal{O}(|F| |G|^2)$
- In step2, the complexity is bound by the number of subforest $\sum_{v \in TopLight(G)} R(F(v), G)$

Time Complexity (An Optimal Decomposition Algorithm)

Analysis

We have established that if $|F| \geq |G|$

$$R(F, G) \leq |G|^2 |F| + \sum_{v \in \text{TopLight}(F)} R(F(v), G)$$

Otherwise

$$R(F, G) \leq |F|^2 |G| + \sum_{w \in \text{TopLight}(G)} R(F, G(w))$$

Time Complexity (An Optimal Decomposition Algorithm)

Analysis

LEMMA 17

$$R(F, G) \leq 4(|F| |G|)^{\frac{3}{2}}$$

Proof. By applying the induction hypothesis on $R(F(v), G)$

$$\begin{aligned} R(F, G) &\leq |G|^2 |F| + \sum_{v \in \text{TopLight}(F)} 4(|F(v)| |G|)^{\frac{3}{2}} \\ &= |G|^2 |F| + 4 |G|^{\frac{3}{2}} \sum_{v \in \text{TopLight}(F)} |F(v)|^{\frac{3}{2}} \\ &\leq |G|^2 |F| + 4 |G|^{\frac{3}{2}} \sum_{v \in \text{TopLight}(F)} |F(v)| \max_{v \in \text{TopLight}(F)} \sqrt{|F(v)|} \\ &\leq |G|^2 |F| + 4 |G|^{\frac{3}{2}} |F| \sqrt{\frac{|F|}{2}} = |G|^2 |F| + \sqrt{8} (|F| |G|)^{\frac{3}{2}} \\ &\leq 4(|F| |G|)^{\frac{3}{2}} \end{aligned}$$

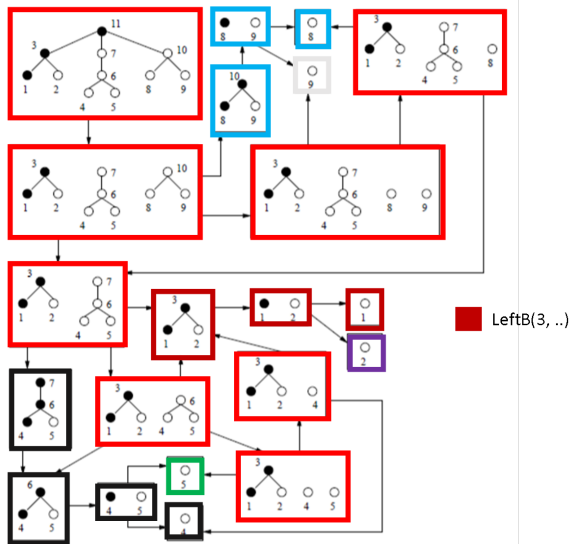
For tree-tree edit distance

$$\delta(l(f), l'(f')) = \min \begin{cases} \delta(f, l'(f')) + \gamma(l \rightarrow \lambda) \\ \delta(l(f), f') + \gamma(\lambda \rightarrow l') \\ \delta(f, f') + \gamma(l \rightarrow l') \end{cases}$$

For forest-forest edit distance

$$\delta(l(f) \circ t, l'(f') \circ t') = \min \begin{cases} \delta(f \circ t, l'(f') \circ t') + \gamma(l \rightarrow \lambda) \\ \delta(l(f) \circ t, f' \circ t') + \gamma(\lambda \rightarrow l') \\ \delta(l(f), l'(f')) + \delta(t, t') \end{cases}$$

Decomposition Example



Definition

- **Free** : the number of sub-problems of tree A compared with tree B.
- **LeftA** : the number of sub-problems created by **right** deletions of tree B (favorite child is the **leftmost** root of tree B), which are compared with all **leftmost** forest of tree A.
- **RightA** : the number of sub-problems created by **left** deletions of tree B (favorite child is the **rightmost** root of tree B), which are compared with all **rightmost** forest of tree A.
- **AllA** : the number of sub-problems after a series of left and right deletions of tree B (favorite child is neither rightmost nor leftmost root of tree B), which are compared with all **special** forest of tree A.

Note

- The definitions of **LeftB**, **RightB**, **AllA** are analogy with **LeftA**, **RightB**, **AllA**.
- The size of each matrix is $|A| * |B|$.

Number of Rightmost, Leftmost and Subforest.

Let A be a tree of size n ,

$$\#right(A) = \sum(|A(i)|, i \in A) - \sum(|A(j)|, j \text{ is a rightmost child})$$

$$\#left(A) = \sum(|A(i)|, i \in A) - \sum(|A(j)|, j \text{ is a leftmost child})$$

$$\#special(A) = \frac{n(n+3)}{2} - \sum_{i \in A} |A(i)|$$

Time Complexity

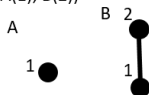
To compute the number of rightmost, leftmost and special forests for each subtree whose each root is a node in tree A and tree B . This can be made in $\mathcal{O}(|A| + |B|)$

Computation of Free

Free(A(1), B(1))

$$\begin{array}{c} A \\ 1 \end{array} \bullet \quad \begin{array}{c} B \\ 1 \end{array} \bullet = 1$$

Free(A(1), B(2))



= min

LeftA(A(1), B(1))



+ #left(A) **favorite child: B1, right decomposition**

RightA(A(1), B(1))



+ #right(A) **favorite child: B1, left decomposition**

Computation of Free

Free(A(1), B(4))

$$\begin{array}{c} \text{A} \quad \text{B} \end{array} \quad \begin{array}{c} 4 \\ \swarrow \downarrow \searrow \\ 1 \quad 2 \quad 3 \end{array} = \begin{array}{c} \text{Free}(A(1), B(2)) \\ 1 \quad 2 \end{array} + \begin{array}{c} \text{Free}(A(1), B(3)) \\ 1 \quad 3 \end{array} + \begin{array}{c} \text{LeftA}(A(1), B(1)) \\ 1 \quad 1 \end{array} + \# \text{left}(A) * (|B| - |B(1)|)$$

favorite child: B1, right decomposition

$$\begin{array}{c} \text{A} \quad \text{B} \end{array} \quad \begin{array}{c} 4 \\ \swarrow \downarrow \searrow \\ 1 \quad 2 \quad 3 \end{array} = \begin{array}{c} \text{Free}(A(1), B(1)) \\ 1 \quad 2 \end{array} + \begin{array}{c} \text{Free}(A(1), B(3)) \\ 1 \quad 3 \end{array} + \begin{array}{c} \text{AllA}(A(1), B(2)) \\ 1 \quad 1 \end{array} + \min \begin{array}{l} \# \text{right}(A) * (1 + |B(1)|) + \# \text{special}(A) * (|B(3)|) \\ \# \text{left}(A) * (1 + |B(3)|) + \# \text{special}(A) * (|B(1)|) \end{array}$$

favorite child: B2, left decomposition

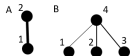
favorite child: B2, right decomposition

$$\begin{array}{c} \text{A} \quad \text{B} \end{array} \quad \begin{array}{c} 4 \\ \swarrow \downarrow \searrow \\ 1 \quad 2 \quad 3 \end{array} = \begin{array}{c} \text{Free}(A(1), B(1)) \\ 1 \quad 2 \end{array} + \begin{array}{c} \text{Free}(A(1), B(2)) \\ 1 \quad 3 \end{array} + \begin{array}{c} \text{RightA}(A(1), B(2)) \\ 1 \quad 1 \end{array} + \# \text{right}(A) * (|B| - |B(3)|)$$

favorite child: B3, left decomposition

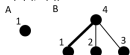
Computation of Free

Free(A(2), B(4))



=

LeftB(A(1), B(4))



#left(B)

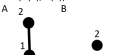
favorite child: A1, right decomposition

+ min

#right(B)

favorite child: A1, left decomposition

Free(A(2), B(2))



+

2

+

3

+

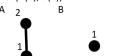
1

+

#left(A) * (|B| - |B(1)|)

favorite child: B1, right decomposition

Free(A(2), B(1))



+

1

+

2

+

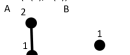
3

+

#right(A) * (|B| - |B(3)|)

favorite child: B3, left decomposition

Free(A(2), B(1))



+

1

+

3

+

2

+

#left(A) * (1 + |B(1)|) + #special(A) * (|B(3)|)

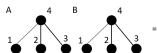
favorite child: B2, right decomposition

+ min

#right(A) * (1 + |B(3)|) + #special(A) * (|B(1)|) favorite child: B2, left decomposition

Computation of Free

Free(A(4), B(4))



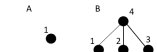
Free(A(2), B(4))



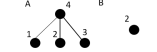
Free(A(1), B(4))



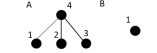
Free(A(1), B(4))



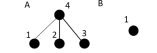
Free(A(4), B(2))



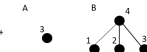
Free(A(4), B(1))



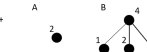
Free(A(4), B(1))



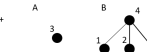
Free(A(3), B(4))



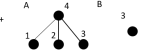
Free(A(2), B(4))



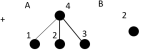
Free(A(3), B(4))



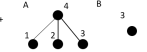
Free(A(4), B(3))



Free(A(4), B(2))



Free(A(4), B(3))



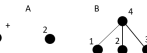
LeftB(A(1), B(4))



RightB(A(3), B(4))



AlIB(A(2), B(4))



LeftA(A(4), B(1))



RightA(A(4), B(3))



AlIA(A(4), B(2))



#left(B) * (|A| - |A(1)|)

favorite child: A1, right decomposition

#right(B) * (|A| - |A(3)|)

favorite child: A3, left decomposition

#left(B) * (1 + |A(3)|) + #special(B) * |A(1)|

favorite child: A2, right decomposition

Or

#right(B) * (1 + |A(1)|) + #special(B) * |A(3)|

favorite child: A2, left decomposition

#left(A) * (|B| - |B(1)|)

favorite child: B1, right decomposition

#right(A) * (|B| - |B(3)|)

favorite child: B3, left decomposition

#left(A) * (1 + |B(3)|) + #special(A) * |B(1)|

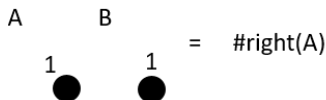
favorite child: B2, right decomposition

#right(A) * (1 + |B(1)|) + #special(A) * |B(3)|

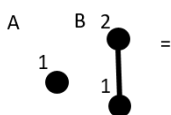
favorite child: B2, left decomposition

Computation of RightA

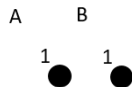
RightA(A(1), B(1))



RightA(A(1), B(2))

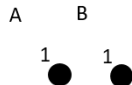


RightA(A(1), B(1))



+ #right(A) **favorite child: B1, left decomposition**

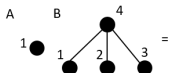
AllA(A(1), B(1))



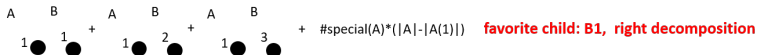
+ #special(A) **favorite child: B1, right decomposition**

Computation of RightA

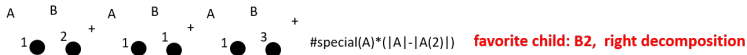
RightA(A(1), B(4))



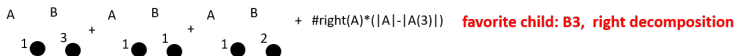
AllA(A(1), B(1)) Free(A(1), B(2)) Free(A(1), B(3))



AllA(A(1), B(2)) Free(A(1), B(1)) Free(A(1), B(3)) #right(A)*(1+|A(1)|) + #special(A)*(|A(3)|) **favorite child: B2, left decomposition**

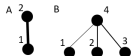


RightA(A(1), B(3)) Free(A(1), B(1)) Free(A(1), B(2))



Computation of RightA

RightA(A(2), B(4))



=

Free(A(2), B(2))



+

Free(A(2), B(3))



+

AllA(A(2), B(1))



+

#special(A) * (|B| - |B(1)|)

favorite child: B1, right decomposition

Free(A(2), B(1))



+

Free(A(2), B(2))



+

RightA(A(2), B(3))



+

#right(A) * (|B| - |B(3)|)

favorite child: B3, left decomposition

Free(A(2), B(1))



+

Free(A(2), B(3))



+

AllA(A(2), B(2))



+

#special(A) * (|B| - |B(2)|)

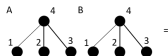
favorite child: B2, right decomposition

+ min

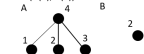
#right(A) * (1 + |B(1)|) + #special(A) * (|B(3)|) **favorite child: B2, left decomposition**

Computation of RightA

RightA(A(4), B(4))



Free(A(4), B(2))



Free(A(4), B(3))



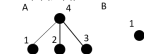
AllA(A(4), B(1))



$\#special(A) * (|B| - |B(1)|)$

favorite child: B1, right decomposition

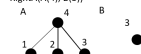
Free(A(4), B(1))



Free(A(4), B(2))



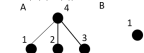
RightA(A(4), B(3))



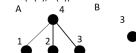
$\#right(A) * (|B| - |B(3)|)$

favorite child: B3, left decomposition

Free(A(4), B(1))



Free(A(4), B(3))



AllA(A(4), B(2))



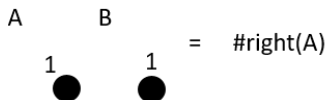
$\#special(A) * (|B| - |B(2)|)$

favorite child: B3, right decomposition

$\#right(A) * (1 + |B(1)|) + \#special(A) * |B(3)|$

Computation of RightA

RightA(A(1), B(1))



RightA(A(1), B(2))

RightA(A(1), B(1))



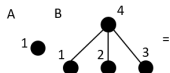
AllA(A(1), B(1))

A B

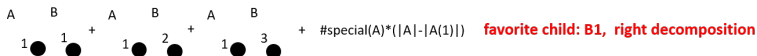


Computation of RightA

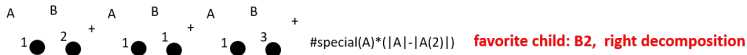
RightA(A(1), B(4))



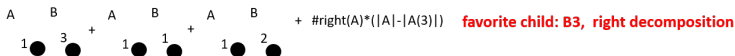
AllA(A(1), B(1)) Free(A(1), B(2)) Free(A(1), B(3))



AllA(A(1), B(2)) Free(A(1), B(1)) Free(A(1), B(3)) #right(A)*(1+|A(1)|) + #special(A)*(|A(3)|) **favorite child: B2, left decomposition**

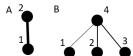


RightA(A(1), B(3)) Free(A(1), B(1)) Free(A(1), B(2))



Computation of RightA

RightA(A(2), B(4))



=

Free(A(2), B(2))



+

Free(A(2), B(3))



+

AllA(A(2), B(1))



+

#special(A) * (|B| - |B(1)|)

favorite child: B1, right decomposition

Free(A(2), B(1))



+

Free(A(2), B(2))



+

RightA(A(2), B(3))



+

#right(A) * (|B| - |B(3)|)

favorite child: B3, left decomposition

Free(A(2), B(1))



+

Free(A(2), B(3))



+

AllA(A(2), B(2))



+

#special(A) * (|B| - |B(2)|)

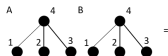
favorite child: B2, right decomposition

+ min

#right(A) * (1 + |B(1)|) + #special(A) * (|B(3)|) **favorite child: B2, left decomposition**

Computation of RightA

RightA(A(4), B(4))



Free(A(4), B(2))



Free(A(4), B(3))



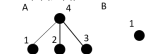
AllA(A(4), B(1))



$\#special(A) * (|B| - |B(1)|)$

favorite child: B1, right decomposition

Free(A(4), B(1))



Free(A(4), B(2))



RightA(A(4), B(3))



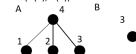
$\#right(A) * (|B| - |B(3)|)$

favorite child: B3, left decomposition

Free(A(4), B(1))



Free(A(4), B(3))



AllA(A(4), B(2))



$\#special(A) * (|B| - |B(2)|)$

favorite child: B3, right decomposition

$\#right(A) * (1 + |B(1)|) + \#special(A) * |B(3)|$

Tree Indexing

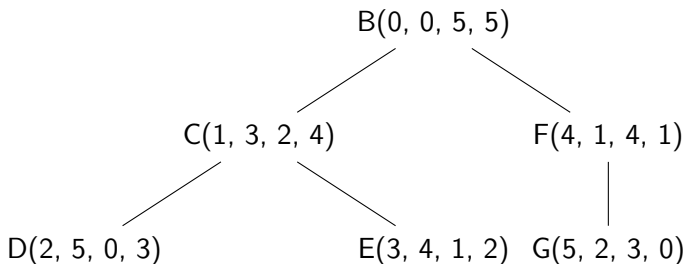
- $preOrder_left_to_right \rightarrow preL$
- $preOrder_right_to_left \rightarrow preR$
- $postOrder_left_to_right \rightarrow postL$
- $postOrder_right_to_left \rightarrow postR$

Note:

- $preL + postR + 1 = treeSize$
- $preR + postL + 1 = treeSize$

Tree Indexing

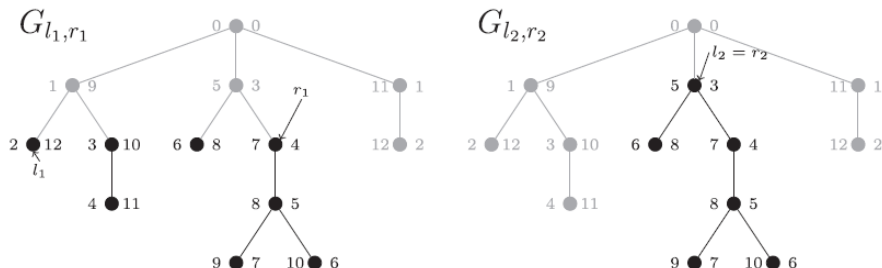
A sequence "**(B(C(D)(E))(F(G)))**" can be built to a tree F.
Consider a tree F(preL, preR, postL, postR)



Note:

- *preOrder_left_to_right* \rightarrow *preL* \rightarrow *lF*
- *preOrder_right_to_left* \rightarrow *preR* \rightarrow *rF*

Tree Indexing

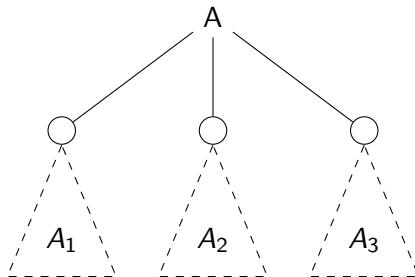


Definition:

$F_{IF, rF}$ is the subforest of F with nodes $N(F_{IF, rF}) = \{IF, rF\} \cup \{x: x \in F, x \text{ succeeds } IF \text{ in left-to-right preorder and } x \text{ succeeds } rF \text{ in right-to-left preorder}\}$ and edges $E(F_{IF, rF}) = \{(v, w) \in E(F): v \in F_{IF, rF} \cap w \in F_{IF, rF}\}$

Computation of Rightmost, Leftmost and Special Forest

Let A be a tree of size n ,



$$\#right(l(A_1, \dots, A_n)) = \sum \#right(A_i) + |A| - |A_n|$$

$$\#left(l(A_1, \dots, A_n)) = \sum \#left(A_i) + |A| - |A_1|$$

$$\#special(A) = \frac{n(n+3)}{2} - \sum_{i \in A} |A(i)|$$

Computation of Rightmost, Leftmost and Special Forest

Let A be a tree of size n ,

$$\#right(l(A_1, \dots, A_n)) = \sum \#right(A_i) + |A| - |A_n|$$

$$\#left(l(A_1, \dots, A_n)) = \sum \#left(A_i) + |A| - |A_1|$$

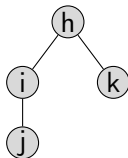
$$\#special(A) = \frac{n(n+3)}{2} - \sum_{i \in A} |A(i)|$$

Note: Each node marks

- the size of subtree rooted at that node
- the sum of the size of subtree of its children
- the number of leftmost forest of subtree rooted at that node
- the number of rightmost forest of subtree rooted at that node
- the number of special forest of subtree rooted at that node

Special Forest Enumeration

Let G be a tree,



Then the **left-to-right** enumeration would be

- $lG = k$

\textcircled{k} ($rG = k$)

- $lG = j$

\textcircled{j} ($rG = j$)

$\textcircled{j} \textcircled{k}$ ($rG = k$)

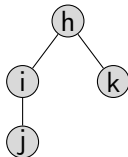
- $lG = i$

\textcircled{i}
 \textcircled{j} ($rG = i$)

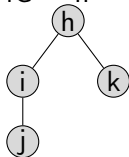
\textcircled{i}
 $\textcircled{j} \textcircled{k}$ ($rG = k$)

Special Forest Enumeration

Let G be a tree,



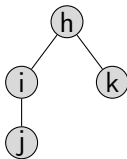
- $IG = h$



$(rG = h)$

Special Forest Enumeration

Let G be a tree,

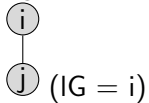


Then the **right-to-left** enumeration would be

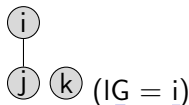
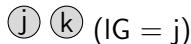
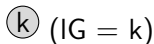
- $rG = j$



- $rG = i$

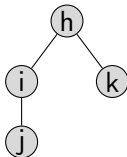


- $rG = k$

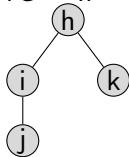


Special Forest Enumeration

Let G be a tree,



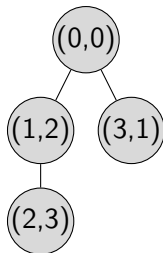
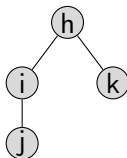
• $rG = h$



$(lG = h)$

Special Forest Enumeration

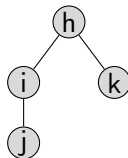
Let G be a tree,



- left-to-right
 - IG enumerate in reverse left to right pre-order.
 - rG enumerate from subset $\text{right}(G, \text{IG}) \cup \text{IG}$ in reverse right to left pre-order.
- right-to-left
 - rG enumerate in reverse right to left pre-order.
 - IG enumerate from subset $\text{left}(G, \text{rG}) \cup \text{rG}$ in reverse left to right pre-order

Rightmost Forest Enumeration - right path decomposition

Let G be a tree,



- $rG = j$
 \textcircled{j} ($IG = j$)

- $rG = i$
 \textcircled{i}
 \textcircled{j} ($IG = i$)

- $rG = k$

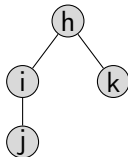
\textcircled{k} ($IG = k$)

$\textcircled{j} \textcircled{k}$ ($IG = j$)

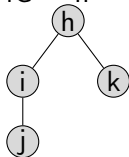
\textcircled{i}
 $\textcircled{j} \textcircled{k}$ ($IG = i$)

Rightmost Forest Enumeration - right path decomposition

Let G be a tree,



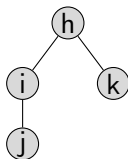
- $lG = h$



$(rG = h)$

Leftmost Forest Enumeration - left path decomposition

Let G be a tree,



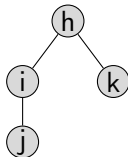
- $lG = k$
 $\overset{j}{\underset{\cdot}{\textcircled{j}}} \text{ (} rG = k \text{)}$

- $lG = j$
 $\overset{i}{\underset{\cdot}{\textcircled{j}}} \text{ (} rG = j \text{)}$

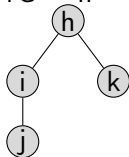
- $lG = i$
 $\overset{i}{\underset{\cdot}{\textcircled{j}}} \text{ (} rG = i \text{)}$ $\overset{i}{\underset{\cdot}{\textcircled{j}}} \textcircled{k} \text{ (} rG = k \text{)}$

Leftmost Forest Enumeration - left path decomposition

Let G be a tree,



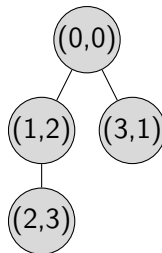
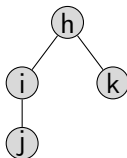
• $rG = h$



$(lG = h)$

Leftmost/Rightmost Forest Enumeration

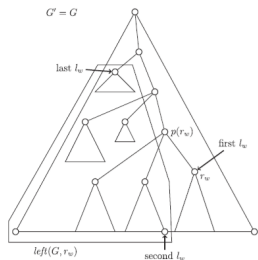
Let G be a tree,



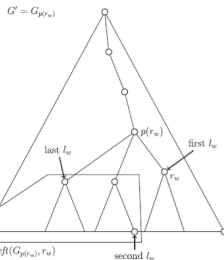
- leftmost forest enumeration
 - IG enumerate in reverse left to right pre-order.
 - rG enumerate from subset $\text{right}(G', IG) \cup IG$ in reverse right to left pre-order.
- rightmost forest enumeration
 - rG enumerate in reverse right to left pre-order.
 - IG enumerate from subset $\text{left}(G', IG) \cup IG$ in reverse left to right pre-order.

Leftmost/Rightmost Forest Enumeration

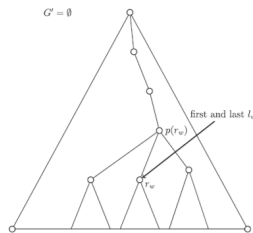
G'



(a) $G' = G$



(b) $G' = G_{p(r_w)}$



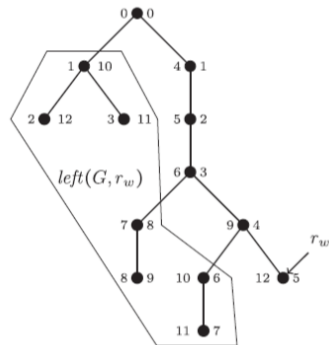
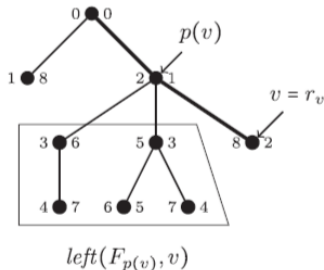
(c) $G' = \emptyset$

- special forest enumeration $G' = G$
- leftmost forest enumeration
 - lG is the leftmost child of its parent $G' = G_p(lG)$
 - otherwise $G' = \emptyset$
- rightmost forest enumeration
 - rG is the rightmost child of its parent $G' = G_p(rG)$
 - otherwise $G' = \emptyset$

Leftmost/Rightmost Forest Enumeration

$\text{left}(G, w)$

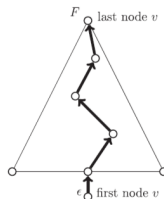
$\text{left}(F_{p(v)}, v)$



Sub Problem Enumeration

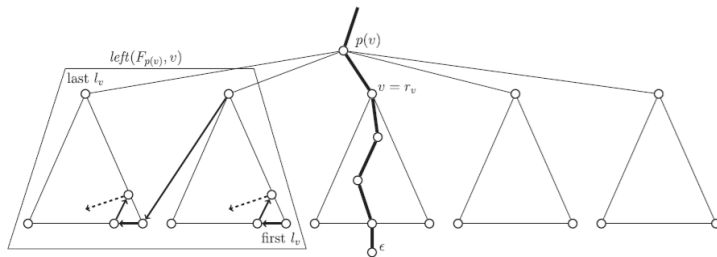
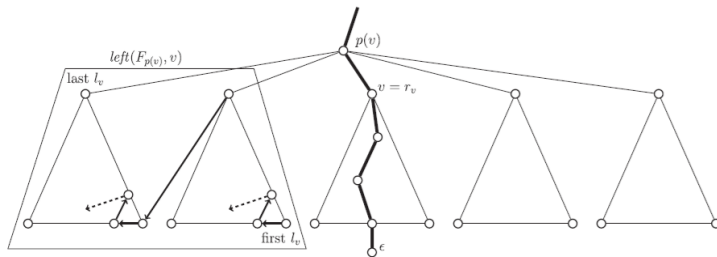
Enumerate quadruples (lF, rF, lG, rG) in loops $(A(BCD)(B'C'D))$.

- Loop A is the outermost loop and iterates bottom-up over the nodes of path.



- Loop C iterates over the nodes to the left of path , in particular, over the nodes $\text{left}(Fp(v), v)$ for a given node v , where v is the index variable of loop A.
- Loop C' is symmetric to loop C and iterates over the nodes $\text{right}(Fp(v), v) \cup p(v)$.

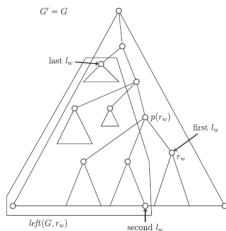
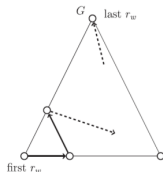
Subproblem Enumeration(Tree F Enumeration)



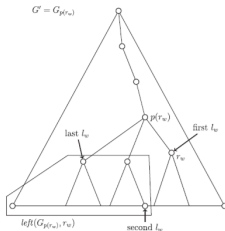
Subproblem Enumeration

Enumerate quadruples (lF, rF, lG, rG) in loops (A(BCD)(B'C'D)).

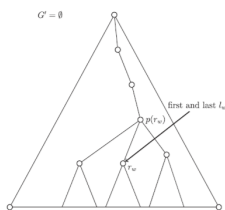
- Loops B and D define subforests of the right-hand tree G.



(a) $G' = G$



(b) $G' = G_{p(r_w)}$

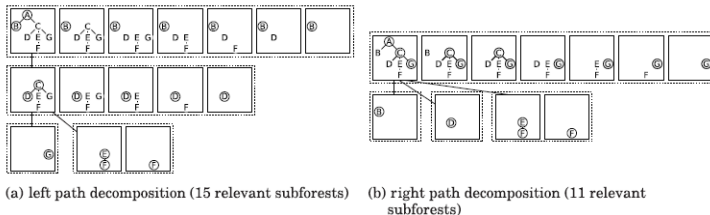


(c) $G' = \emptyset$

- Time complexity in average $\mathcal{O}(n^3)$:
 - The path decomposition in F create n sub forests.
 - The all decomposition in G create n^2 sub forests.
- The naive implement leads to the space complexity in average $\mathcal{O}(n^3)$, but can be reduced to $\mathcal{O}(n^2)$

The intermediate subforest enumeration

The intermediate subforest enumeration with respect to $F(V_{p-1})$ and $F(V_p)$ is the sequence of forests $F(V_{p-1}) = F_0, F_1 \cdots F_5 = F(V_p)$



The intermediate subforest enumeration

$$\delta_L(F_{IF,rF}, G_{IG,rG}) = \min \begin{cases} \delta(F_{IF,rF} - IF, G_{IG,rG}) + \gamma(IF \rightarrow \lambda) \\ \delta(F_{IF,rF}, G_{IG,rG} - IG) + \gamma(\lambda \rightarrow IG) \\ \delta(F_{IF,rF} - F_{IF}, G_{IG,rG} - F_{IG}) \\ \quad + \delta(F_{IF} - IF, G_{IG} - IG) + \gamma(IF \rightarrow IG) \end{cases}$$

$\mathcal{O}(n^2)$ Space Complexity Implement

Four Memorization tables: D, S, T, and Q

- D(of size $|F| |G|$) stores distance values which are needed across different calls of the single-path function. Subproblems in D are **pairs of subtrees without their root nodes**.
- S(of size $|F| |G|$) stores all subproblems computed in the two innermost loops C(IF in left decompose) and D(IG in left decompose) or C'(rF in right decompose) and D'(rG in right decompose) for specific nodes in the outer loops. → maintain in each round in B and B'.
- T(of size $|G|^2$) stores distances between **a specific subforest in F** and **each relevant subforest of G**. The subforest is either a subtree rooted on the path or a subforest $F_{l'_v, v}$ where l'_v is the leftmost child of $p(v)$.
- Q(of size $|F|$) stores distances between **each subforest of F defined in loop C**(path decomposition) and **a specific subforest** in G of the form $G_{p(r_G)-p(r_G)}$, that **is a subtree in G without root node**. used for the next C loop.

$\mathcal{O}(n^2)$ Space Complexity Implement

Tables indexing:

For $\delta_L(F_{IF,rF}, G_{IG,rG})$

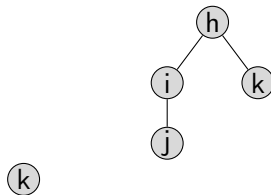
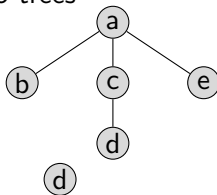
- D: $D[\text{preR_to_preL}(rF), \text{preR_to_preL}[rG]]$
- S: $S[IF, IG]$
- T: $T[IG, rG]$
- Q: $Q[IF]$

For $\delta_R(F_{IF,rF}, G_{IG,rG})$

- D: $D[IF, rF]$
- S: $S[rF, rG]$
- T: $T[IG, rG]$
- Q: $Q[rF]$

$O(n^2)$ Space Complexity Implement - Example

Let F and G be two trees



$T[$