An Efficient Algorithm for Computing Non-overlapping Inversion and Transposition Distance

Toan Thang Ta, Cheng-Yao Lin and Chin Lung Lu

Department of Computer Science National Tsing Hua University, Hsinchu 30013, Taiwan {toanthanghy, begoingto0830}@gmail.com, cllu@cs.nthu.edu.tw

Abstract

Given two strings of the same length n, the non-overlapping inversion and transposition distance (also called mutation distance) between them is defined as the minimum number of non-overlapping inversion and transposition operations used to transform one string into the other. In this study, we present an $O(n^3)$ time and $O(n^2)$ space algorithm to compute the mutation distance of two input strings.

1 Introduction

The dissimilarity of two strings is usually measured by the so-called edit distance, which is defined as the minimum number of edit operations necessary to convert one string into the other. The commonly used edit operations are character insertions, deletions and substitutions. In biological application, the aforementioned edit operations correspond to point mutations of DNA sequences (i.e., mutations at the level of individual nucleotides). From evolutionary point of view, however, DNA sequences may evolve by large-scale mutations (also called rearrangements, i.e., mutations at the level of sequence fragments) [8], such as inversions (replacing a fragment of DNA sequence by its reverse complement) and transpositions (i.e., moving a fragment of DNA sequence from one location to another or, equivalently, exchanging two adjacent and non-overlapping fragments on DNA sequence). Note that a large-scale mutation that replaces a fragment of DNA sequence only by its reverse (without complement) is called a reversal. Based large-scale on mutation operations, the dissimilarity (or mutation distance) between two strings can be defined to be the minimum number of large-scale mutation operations used to transform one string to the other. Cantone et al. [2] introduced an O(nm) time and $O(m^2)$ space algorithm to solve an approximate string matching problem with non-overlapping reversals, which is to find all locations of a given text that match a

given pattern with non-overlapping reversals, where n is the length of the text and m is the length of the pattern. In this problem, two equal-length strings are said to have a match with non-overlapping reversals if one string can be transformed into the other using any finite sequence of non-overlapping reversals. It should be noted that the number of the used non-overlapping reversals in the algorithm proposed by Cantone et al. [2] is not required to be less than or equal to a fixed non-negative integer. In [2], Cantone et al. also presented another algorithm whose average-case time complexity is O(n). Cantone et al. [3] studied the same problem by considering both non-overlapping reversals and transpositions, where they called transpositions as translocations and the lengths of two exchanged adjacent fragments are constrained to be equivalent. They finally designed an algorithm to solve this problem in $O(nm^2)$ time and $O(m^2)$ space. For the above problem, Grabowski et at. [6] gave another algorithm whose worst-case time and space are $O(nm^2)$ and O(m), respectively. Moreover, they proved that their algorithm has an O(n) average time complexity. Recently, Huang et al. [7] studied the above approximate string problem under non-overlapping matching reversals by further restricting the number of the used reversals not to exceed a given positive integer k. They proposed a dynamic programming algorithm to solve this problem in $O(nm^2)$ time and $O(m^2)$ space.

In this work, we are interested in the computation of the mutation distance between two strings of the same length under non-overlapping inversions and transpositions (i.e., non-overlapping inversion and transposition distance), where the lengths of two adjacent and non-overlapping fragments exchanged by a transposition can be different. For this problem, we devise an algorithm whose time and space complexities are $O(n^3)$ and $O(n^2)$, respectively, where n is the length of two input strings. The rest of the paper is organized as follows. In Section 2, we provide some notations that are helpful when we present our algorithm later. Next, we develop the main algorithm and also analyze its time and space complexities in Section 3. Finally, we have a brief conclusion in Section 4.

2 **Preliminaries**

Let x be a string of length n over a finite alphabet Σ . A character at position *i* of *x* is represented with x_i , where $1 \le i \le n$. A substring of x from position i to j is indicated as $x_{i,j}$, i.e., $x_{i,j} = x_i x_{i+1} \dots x_j$, for $1 \le i \le j \le n$. In biological sequences, $\Sigma = \{A, C, G, T\}$, in which A - T and C - G are considered as complementary base pairs. We use $\theta(x)$ to denote an inversion operation acting on a string x, resulting in a reverse and complement of x. For example, $\theta(A) = T$, $\theta(T) = A$, $\theta(G) = C$, $\theta(C) = G$ and $\theta(CGA) = TCG$. In addition, we utilize $\tau(uv) = vu$ to represent a transposition operation to exchange two non-empty strings uand v. Note that the lengths of u and v are required to be identical in some previous works [3, 5, 6], but they may be different in this study. For convenience, we call θ and τ as mutation operations. We also let $\theta_{i,i}(x) = \theta(x_{i,i})$ for $1 \leq 1$ $i \leq j \leq n$ and $\tau_{i,j,k}(x) = x_{k,j}x_{i,k-1}$ for $1 \leq i < i$ $k \le j \le n$, where [i, j] is called a *mutation range* for $\theta_{i,i}$ and $\tau_{i,i,k}$ and k is called as a *cut point* in $\tau_{i,j,k}$.

For an integer $1 \le t \le n$, we say that a mutation operation $\theta_{i,j}$ or $\tau_{i,j,k}$ covers t if $i \le t \le j$. Given two mutation operations, they are *non-overlapping* if the intersection of their mutation ranges are empty. In this study, we are only interested in non-overlapping mutation operations. Consider a set Θ of non-overlapping mutation operations and a string x, let $\Theta(x)$ be the resulting string after consecutively applying mutation operations in Θ on x. For example, suppose that $\Theta = \{\tau_{1,3,2}, \theta_{5,5}\}$ and x = TAGAC. Then we have $\Theta(x) = AGTAG$.

Definition 1. (Non-overlapping inversion and transposition distance) Given two strings x and y of the same length, the non-overlapping inversion and transposition distance (simply called *mutation distance*) between x and y, denoted by md(x, y), is defined as the minimum number of non-overlapping inversion and transposition operations used to transform x into y. If there does not exist any set of non-overlapping

mutations that converts x into y, then md(x, y) is infinite. Formally,

md(x, y) =

 $\begin{cases}
\min\{|\Theta| : \Theta(x) = y\} & \text{if } \exists \Theta \text{ such that } \Theta(x) = y \\
\infty & \text{otherwise}
\end{cases}$

For example, consider x = TAGAC and y = TAACG. Then there are only two sets of non-overlapping mutation operations $\Theta_1 = \{\theta_{1,2}, \tau_{3,5,4}\}$ and $\Theta_2 = \{\tau_{3,5,4}\}$ such that $\Theta_1(x) = y$ and $\Theta_2(x) = y$. Thus, $md(x, y) = |\Theta_2| = 1$.

3 The algorithm

transposition (respectively, Basically, а inversion) operation acting on a string x can be considered as a permutation of characters in x(respectively, complement of x). From this view point, a mutation operation actually comprises several sub-operations, called as mutation fragments, each of which is denoted either by a tuple (i, j, x_j) or $(i, j, \theta(x_j))$. The mutation fragment (i, j, x_i) (respectively, $(i, j, \theta(x_i))$) means that x_i (respectively, complement of x_i) is moved into the position *i* in the resulting string obtained when applying the mutation operation on x. For convenience, we arrange all possible mutation fragments in a three-dimensional table $M_x[2,n,n]$, which is called *mutation table* of x, as follows.

- $M_x[1, i, j] = (i, j, \theta(x_i))$ for i, j = 1, 2, ..., n.
- $M_x[2, i, j] = (i, j, x_i)$ for i, j = 1, 2, ..., n.

For example, let x = TAGAC. Then its mutation table is shown in Figure 1.

When an inversion operation $\theta_{i,j}$ applies on a string x, we can decompose it into (j - i + 1)mutation fragments $\mathcal{F}(\theta_{i,i}, x, t)$ for $i \leq t \leq j$, where $\mathcal{F}(\theta_{i,j}, x, t) = (t, i+j-t, \theta(x_{i+j-t}))$. Similarly, a transposition operation $\tau_{i,i,k}$ can be also decomposed into (j - i + 1) mutation fragments $\mathcal{F}(\tau_{i,j,k}, x, t)$ for $i \leq t \leq j$, where if $i \le t \le i + j - k$, then $\mathcal{F}(\tau_{i,j,k}, x, t) = (t, k + j)$ $t - i, x_{k+t-i}$; otherwise (i.e., if $i + j - k < t \le j$ $\mathcal{F}(\tau_{i,j,k}, x, t) = (t, t - j + k - j)$ j), then 1, $x_{t-i+k-1}$). For the sake of succinctness, let $\tau_{i,j,k}(x,1) = \{(t,k+t-i,x_{k+t-i}) : i \le t \le i + i\}$ j-kand $\tau_{i,j,k}(x,2) = \{(t,t-j+k 1, x_{t-j+k-1}) : i+j-k < t \le j \}.$

$M_x[1, i, j]$	i = 1	2	3	4	5
j = 1	(1,1,A)	(2,1,A)	(3,1, <i>A</i>)	(4, 1, A)	(5,1, <i>A</i>)
2	(1,2,T)	(2,2,T)	(3,2, <i>T</i>)	(4, 2, T)	(5,2, <i>T</i>)
3	(1,3, <i>C</i>)	(2,3, <i>C</i>)	(3,3, <i>C</i>)	(4,3, <i>C</i>)	(5,3, <i>C</i>)
4	(1, 4, T)	(2,4 <i>,T</i>)	(3,4 <i>,T</i>)	(4, 4, T)	(5,4 <i>,T</i>)
5	(1,5,G)	(2,5,G)	(3,5,G)	(4,5,G)	(5,5,G)

$M_x[2, i, j]$	i = 1	2	3	4	5
j = 1	(1,1,T)	(2,1,T)	(3,1,T)	(4, 1, T)	(5, 1, T)
2	(1,2,A)	(2,2, <i>A</i>)	(3,2, <i>A</i>)	(4, 2, A)	(5,2,A)
3	(1,3,G)	(2,3,G)	(3,3,G)	(4,3,G)	(5,3,G)
4	(1,4, <i>A</i>)	(2,4, <i>A</i>)	(3,4, <i>A</i>)	(4,4, <i>A</i>)	(5,4,A)
5	(1,5, <i>C</i>)	(2,5, <i>C</i>)	(3,5, <i>C</i>)	(4,5, <i>C</i>)	(5,5, <i>C</i>)

Figure 1. A mutation table M_x for a given string x = TAGAC, where the column is indexed by *i* and the row by *j*. Shaded cells respectively represent the inversion $\theta_{1,3}(x)$ on $M_x[1, i, j]$ and the transposition operation $\tau_{1,5,4}(x)$ on $M_x[2, i, j]$.

The aforementioned decomposition can be extended to apply on a set Θ of non-overlapping mutation operations. That is, when Θ acts on a string x, it can be decomposed into a sequence $\mathcal{F}(\Theta, x) = \langle F_1, F_2, \dots, F_n \rangle$ of mutation fragments by the following formula: For $t = 1, 2, \dots, n$,

$$F_t = \begin{cases} \mathcal{F}(\theta_{i,j}, x, t) & \text{if } \exists \theta_{i,j} \in \Theta \text{ that covers } t \\ \mathcal{F}(\tau_{i,j,k}, x, t) & \text{if } \exists \tau_{i,j,k} \in \Theta \text{ that covers } t \\ (t, t, x_t) & \text{otherwise} \end{cases}$$

For instance, given $\Theta = \{\tau_{1,2,2}, \theta_{4,5}\}$ and x = TAGAC, we have $\mathcal{F}(\Theta, x) = \langle (1,2,A), (2,1,T), (3,3,G), (4,5,G), (5,4,T) \rangle$.

Observation 1. Given a string x and its mutation table M_x , the result of an inversion $\theta_{i,i}(x)$ can be obtained by concatenating (j - i + 1)mutation fragments starting at $M_{x}[1, i, j]$ and continuing to move in the anti-diagonal direction to $M_x[1, j, i]$. Moreover, the result of a transposition $\tau_{i,i,k}(x)$ is obtained by concatenating the mutation fragments on the following two paths, one starting at $M_r[2, i, k]$ and continuing to move in the diagonal direction to $M_{x}[2, i + j - k, j]$ and the other beginning at $M_{x}[2, i+j-k+1, i]$ and also continuing to move in the diagonal direction to $M_x[2, j, k-1]$.

For a mutation fragment $F = (i, j, \sigma)$, we say that F yields the character σ , denoted by $\ell(F) = \sigma$. If a sequence $S = \langle F_1, F_2, \dots, F_m \rangle$ consists of *m* mutation fragments (*m* is also considered as the *length* of *S*) with $\ell(F_i) = \sigma_i$ for $1 \le i \le m$, then we say that *S* yields a string $\sigma_1 \sigma_2 \dots \sigma_m$, which is further written as $\ell(S) =$ $\sigma_1 \sigma_2 \dots \sigma_m$. A subsequence *P* containing first *t* elements of *S*, i.e., $P = \langle F_1, F_2, \dots, F_t \rangle$, is called a *prefix* of *S* and denoted by $P \sqsubset S$, where $1 \le t \le m$.

Definition 2. (Agreed sequence) A sequence $S = \langle F_1, F_2, ..., F_m \rangle$ of *m* mutation fragments is an *agreed sequence* if there exists a set Θ of non-overlapping mutation operations such that $S \sqsubset \mathcal{F}(\Theta, x)$, where $m \le n$.

Given a mutation operation $\theta_{i,j}$ or $\tau_{i,j,k}$, we call *i* and *j* as the *left end point* and *right end point* of the mutation operation, respectively. We also say that this mutation operation (i.e., $\theta_{i,j}$ or $\tau_{i,j,k}$) *intersects* with a range [a,b] if the intersection of [i,j] and [a,b] is non-empty. The number of mutation operations in Θ that intersects with the range [1,m] where $m \le n$ is denoted as $\rho(\Theta,m)$. For example, if $\Theta = \{\theta_{1,2}, \tau_{4,6,5}, \theta_{7,9}\}$, then $\rho(\Theta, 4) = 2$. Suppose that a mutation fragment $F = (i,j,\sigma)$ is in $\mathcal{F}(\Theta,x)$ for some set Θ of non-overlapping mutation operations. Then *F* is said to be *covered* by a mutation operation $\theta_{l,r}$ or $\tau_{l,r,k}$ in Θ if $F = \mathcal{F}(\theta_{l,r}, x, i)$ or $F = \mathcal{F}(\tau_{l,r,k}, x, i)$. If *F* is not

covered by any mutation operation in Θ , that is $F = (i, i, x_i)$, then we can consider that F is still covered by a *virtual mutation operation*. We use $\mathcal{L}(\Theta, x, F)$ and $\mathcal{R}(\Theta, x, F)$ to denote the left and right end points of a mutation operation in Θ respectively that covers F. Formally,

•
$$\mathcal{L}(\Theta, x, F) =$$

$$\begin{cases}
l & \text{if } \exists \theta_{l,r} \text{ or } \tau_{l,r,k} \in \Theta \text{ s. t. } F = \mathcal{F}(\theta_{l,r}, x, i) \\
or & F = \mathcal{F}(\tau_{l,r,k}, x, i) \\
i & \text{otherwise}
\end{cases}$$

•
$$\mathcal{R}(\Theta, x, F) =$$

$$\begin{cases} \text{if } \exists \theta_{l,r} \text{ or } \tau_{l,r,k} \in \Theta \text{ s. t. } F = \mathcal{F}(\theta_{l,r}, x, i) \\ \text{or } F = \mathcal{F}(\tau_{l,r,k}, x, i) \\ i \text{ otherwise} \end{cases}$$

For example, suppose that $\Theta = \{\tau_{1,3,2}, \theta_{4,5}\}, x = TAGAC$ and $\mathcal{F}(\Theta, x) = \langle (1,2,A), (2,3,G), (3,1,T), (4,5,G), (5,4,T) \rangle$. For F = (2,3,G), we have $\mathcal{L}(\Theta, x, F) = 1$ and $\mathcal{R}(\Theta, x, F) = 3$.

For an integer *i*, let $A_x^i = \{(i + 1, t, \theta(x_t)) : i + 1 < t \le n\}$ and $B_x^i = \{(i + 1, t, x_t) : i + 1 < t \le n\}$, which are sets of mutation fragments covered by inversion and transposition operations respectively acting on *x* whose left end points are all *i* + 1, where $0 \le i < n$.

Definition 3. Let $S = \langle F_1, F_2, ..., F_m \rangle$ be an agreed sequence of m mutation fragments of x, where $m \leq n$, and Θ be a set of non-overlapping mutation operations such that $S \sqsubset \mathcal{F}(\Theta, x)$. Sequence S is called a complete sequence if $\mathcal{R}(\Theta, x, F_m) = m$.

Observation 2. Suppose that $S = \langle F_1, F_2, ..., F_m \rangle$ is a complete sequence, where m < n. Then $S' = \langle F_1, F_2, ..., F_m, F_{m+1} \rangle$ is an agreed sequence if $F_{m+1} \in A_x^m \cup M_x[1, m+1, m+1] \cup B_x^m \cup M_x[2, m+1, m+1]$.

Definition 4. (Legal sequence) An agreed sequence $S = \langle F_1, F_2, ..., F_m \rangle$ of *m* mutation fragments of *x* is *legal* if $\ell(S) = y_{1,m}$, where $m \leq n$.

Observation 3. Assume that $S = \langle F_1, F_2, ..., F_m \rangle$ is a legal sequence. Then $S' = \langle F_1, F_2, ..., F_{m-1} \rangle$ is also a legal sequence and $\ell(F_m) = y_m$, where $1 < m \le n$.

The main idea of our algorithm is as follows. First, the algorithm constructs all legal sequences of length 1 by examining all mutation fragments in $A_x^0 \cup M_x[1,1,1] \cup B_x^0 \cup M_x[2,1,1]$. A legal sequence is then created if the mutation fragment yields y_1 . Next, for each $1 \le i < n$, the algorithm produces all possible legal sequences of length i + 1, which can be obtained from all previously created legal sequences of length i. At the end, if there exists a legal sequence of length n, then the algorithm returns the minimum number of the used mutation operations among all legal sequences of length n. Otherwise, the algorithm returns infinity.

To distinguish all the legal sequences of length i, we define three sets of *extended* mutation fragments I_x^i , V_x^i and T_x^i of x as follows:

- $I_x^i = \{ (g, d, (i, j, \theta(x_j))) \}$ for $0 \le g \le n 1$, $0 \le d \le n$ and j = 1, 2, ..., n.
- $V_x^i = \{(0, d, (i, i, x_i))\}$ for $0 \le d \le n$.

• $T_x^i = \{(s, d, 1, (i, j, x_j))\} \cup \{(g, d, 2, (i, j, x_j))\}$ for $1 \le s \le n - 1, 0 \le g \le n - 1, 0 \le d \le n$ and j = 1, 2, ..., n.

Actually, each extended mutation fragment contains the last mutation fragment F of some agreed sequence S of length i with extra guiding information g, d and s. The value of g is the number of steps that are required to move towards the anti-diagonal or diagonal direction to complete a mutation operation and d indicates the number of the currently used mutation operations in S. If the last mutation fragment F is covered by a transposition operation, say m_x , and $F \in m_x(x, 1)$, then s denotes the left end point of m_x . An extended mutation fragment is further called an *ending fragment* (respectively, *continuing fragment*) if the value of its first component is zero (respectively, non-zero).

Next, we represent each legal sequence by an extended fragment as follows. Given a legal sequence $S = \langle F_1, F_2, ..., F_i \rangle$ that consists of *i* mutation fragments of *x*, let Θ be any set of non-overlapping mutation operations such that $S \sqsubset \mathcal{F}(\Theta, x)$. We further let m_x be the mutation operation in Θ that covers F_i . Note that m_x may be a virtual mutation operation. We create an extended fragment *e* by the following cases (each case is also called as a *type* of extended mutation fragment):

Case 1 (type *I*): m_x is an inversion operation. Then $e = (\mathcal{R}(\Theta, x, F_i) - i, \rho(\Theta, i), F_i).$

Case 2 (type V): m_x is a virtual mutation operation. Then $e = (0, \rho(\Theta, i), F_i)$.

Case 3 (type *F*): m_x is a transposition operation and $F_i \in m_x(x, 1)$. Then $e = (\mathcal{L}(\Theta, x, F_i), \rho(\Theta, i), 1, F_i)$.

Case 4 (type S): m_x is a transposition operation and $F_i \in m_x(x, 2)$. Then $e = (\mathcal{R}(\Theta, x, F_i) - i, \rho(\Theta, i), 2, F_i)$.

Now, we briefly describe how to produce all possible legal sequences of length i + 1 from all previously created legal sequences of length i. Let L_i be the set of all extended mutation fragments at column i, i.e., each element in L_i corresponds to some legal sequence of length i. Suppose that L_i

is already known. Below, we construct L_{i+1} by examining all extended mutation fragments in L_i . For each extended mutation fragment $e \in L_i$, we identify all candidate mutation fragments at column i + 1 of the mutation tables M_x based on the guiding information contained in e. Let $N_x^i(e)$ be set of the candidate mutation fragments of x at column i + 1. Then we compute $N_x^i(e)$ as follows:

Case 1: *e* is a continuing fragment. Let $F_i = (i, j, \sigma)$ be the mutation fragment contained in *e*. Then we consider three sub-cases to compute $N_x^i(e)$: (1) Suppose that *e* is of type *I*. Then $N_x^i(e) = \{(i + 1, j - 1, \theta(x_{j-1}))\}$. (2) Suppose that *e* is of type *F*, i.e., $e = (s, d, 1, F_i)$. If j < n, then $N_x^i(e) = \{(i + 1, j + 1, x_{j+1}), (i + 1, s, x_s)\}$; otherwise (i.e., j = n), $N_x^i(e) = \{(i + 1, s, x_s)\}$. (3) Suppose that *e* is of type *S*, i.e., $e = (g, d, 2, F_i)$. Then $N_x^i(e) = \{(i + 1, j + 1, x_{j+1}), (i + 1, x_{j+1})\}$.

Case 2: *e* is an ending fragment. We set $N_x^i(e) = A_x^i \cup M_x[1, i+1, i+1] \cup B_x^i \cup M_x[2, i+1, i+1].$

For any $F_{i+1} \in N_x^i(e)$ with $\ell(F_{i+1}) = y_{i+1}$, an extended mutation fragment e' is created to contain F_{i+1} . As to the values of the guiding information in e', they can easily be obtained from those in e based on Observations 1 and 2. Finally, we add e' into L_{i+1} .

Lemma 1. Let P_1 and P_2 be two complete sequences of length m of x with $\ell(P_1) = \ell(P_2)$. Moreover, let S_1 be an agreed sequence of x with $P_1 \sqsubset S_1$ and S_2 be the resulting sequence obtained by replacing P_1 with P_2 in S_1 . Then, S_2 is also an agreed sequence of x and $\ell(S_1) =$ $\ell(S_2)$.

Proof. Since S_1 and P_2 are agreed sequences of x, there exist sets of non-overlapping mutation operations Θ_1 and Θ_2 such that $S_1 \sqsubset \mathcal{F}(\Theta_1, x)$ and $P_2 \sqsubset \mathcal{F}(\Theta_2, x)$. Moreover, $P_1 \sqsubset \mathcal{F}(\Theta_1, x)$ since $P_1 \sqsubset S_1$. We construct a set Θ_3 of non-overlapping mutation operations from $\boldsymbol{\Theta}_1$ as follows. First, let Θ_3 be the resulting set obtained by removing all mutation operations in Θ_1 that intersect with the range [1, m]. Note that the ranges of these removed mutation operations are completely contained in the range [1, m] because P_1 is a complete sequence of length *m*. Next, we insert those mutation operations in Θ_2 whose ranges entirely lie in the range [1, m] into Θ_3 . It can be verified that Θ_3 is a set of non-overlapping mutation operations and as a result, $S_2 \sqsubset \mathcal{F}(\Theta_3, x)$ and $\ell(S_1) = \ell(S_2)$.

Given an extended mutation fragment e, we use d(e) to indicate the value of its second component (i.e., the number of the used mutation operations)

for convenience. Based on Lemma 1, we have the following corollary.

Corollary 1. Let e_1 and e_2 be two ending fragments in L_i . If $d(e_1) < d(e_2)$, then we can safely discard e_2 from L_i when computing the mutation distance between x and y.

Lemma 2. Let e_1 and e_2 be two continuing fragments e_1 and e_2 of the same type S in L_i with $e_1 = (g, d_1, 2, (i, j, x_j))$ and $e_2 = (g, d_2, 2, (i, j, x_j))$. If $d(e_1) < d(e_2)$ (i.e., $d_1 < d_2$), then we can safely remove e_2 from L_i without affecting the computation of the mutation distance between x and y.

Proof. Since two continuing fragments e_1 and e_2 of type *S* contain the same mutation fragment and guiding information *g*, the right end point of the mutation operation covering (i, j, x_j) in e_1 is equal to that of the mutation operation covering (i, j, x_j) in e_2 . Denote by *t* the above right end point. Therefore, the candidate mutation fragments from column *i* to column *t* on the mutation table M_x are exactly equivalent. Then we can safely discard e_2 from L_i if $d_1 < d_2$ according to Corollary 1.

Based on Corollary 1, we can also verify that there does not exist two continuing fragments of same type I or F in L_i such that the values of all their components, excluding the second components (i.e., the number of used mutation operations), are equal. Now, the details of our algorithm for computing non-overlapping inversion and transposition distance between two strings is described in Algorithm 1.

Algorithm Computing non-overlapping inversion and transposition distance between two strings **Input:** Two strings x and y of the same length n. **Output**: Mutation distance md(x, y). 1) Let i = 0, d = 0 and $L_1 = \emptyset$; 2) Construct L_1 by calling EMF_Creation(i,d); 3) if $L_1 = \emptyset$ then return ∞ ; 4) for i=1 to n-1 do $L_{i+1} = \text{EMF}_\text{Extension}(L_i);$ 5) if $L_{i+1} = \emptyset$ then return ∞ ; 6) end for 7) 8) return min{ $d(e) : e \in L_n$ }; **Procedure 1**: EMF Creation(*i*, *d*)

Input: Column i and current number of mutation operations d.

Output: Extended mutation fragments to be added to L_{i+1} when there is an ending fragment at column i. for each $(i+1,t,\theta(x_t)) \in A_x^i$ do 1) 2) if $\theta(x_t) = y_{i+1}$ then Add (t - i - 1, d + 1, (i + 1))3) $1, t, \theta(x_t)$) of type *I* to L_{i+1} ; end if 4) 5) end for 6) if $\theta(x_{i+1}) = y_{i+1}$ then 7) Add ending fragment (0, d + $1, (i+1, i+1, \theta(x_{i+1})))$ of type *I* to L_{i+1} based on Corollary 1; end if 8) 9) for each $(i + 1, t, x_t) \in B_x^i$ do 10) if $x_t = y_{i+1}$ then 11) Add $(i + 1, d + 1, 1, (i + 1, t, x_t))$ of type F to L_{i+1} ; end if 12) 13) end for 14) if $x_{i+1} = y_{i+1}$ then Add ending fragment (0, d, (i +15) $1, i + 1, x_{i+1})$ of type V to L_{i+1} based on Corollary 1; 16) end if **Procedure 2:** EMF Extension (L_i) **Input**: L_i . Output: L_{i+1} . 1) Let E_i be a set of ending fragments at column i; $L_{i+1} \leftarrow \emptyset; E_i \leftarrow \emptyset;$ 2) 3) for each extended mutation fragment e in L_i **do case** 1: e is of type F, let e =4) $(s, d, 1, (i, j, x_i)).$ if j < n and $x_{j+1} = y_{i+1}$ then 5) 6) Add $(s, d, 1, (i + 1, j + 1, x_{j+1}))$ of type F to L_{i+1} ; 7) end if 8) if $x_s = y_{i+1}$ then Add (j - i - 1, d, 2, (i + 1))9) $(1, s, x_s))$ of type S to L_{i+1} , based on Lemma 2; 10) end if 11) **case** 2: e is of type S, let e = $(g, d, 2, (i, j, x_i)).$ if g > 0 then 12) 13) if $x_{j+1} = y_{i+1}$ then Add (g - 1, d, 2, (i + 1, j +14) $(1, x_{i+1}))$ of type S to L_{i+1} based on Lemma 2; 15) end if 16) //g = 0else 17) Add e to E_i based on Corollary 1; end if 18)

19) **case** 3: *e* is of type *I*, let $e = (g, d, (i, j, \theta(x_i))).$ if g > 0 then 20) 21) if $\theta(x_{j-1}) = y_{i+1}$ then Add (g - 1, d, (i + 1, j -22) $1, \theta(x_{i-1}))$ of type I to L_{i+1} ; 23) end if //g = 024) else 25) Add e to E_i based on Corollary 1; 26) end if 27) case 4: e is of type V. Add e to E_i based on 28) Corollary 1; 29) end for 30) if $E_i \neq \emptyset$ then $d = \min\{d(e) : e \in E_i\};$ 31) 32) EMF Creation(i, d); 33) end if 34) return L_{i+1} ;

Theorem 1. Algorithm 1 computes mutation distance of two input strings in $O(n^3)$ time and $O(n^2)$ space.

Proof. The correctness of Algorithm 1 can be verified according to the discussion we mentioned before. Below, we analyze the complexities of Algorithm 1. For convenience, we use |H| to denote the number of extended mutation fragments of type H in L_i . By Corollary 1 and Lemma 2, the numbers of extended mutation fragments of type I, F and S in L_i are $O(n^2)$. Therefore, Procedure 2 can be done in O(|I| + $|S| + |F| = O(n^2)$ time. It is not hard to see that Procedure 1 costs O(n) time in the worst case. Moreover, Procedure 2 is called at most O(n)times in Algorithm 1. As a result, the total time complexity of Algorithm 1 is $O(n^3)$. In addition, the space complexity of Algorithm 1 is O(|I| + $|F| + |S|) = O(n^2).$

4 Conclusions

In this work, we studied the computation of non-overlapping inversion and transposition distance between two strings, which can be useful applications especially in computational biology (e.g., evolutionary tree reconstruction of biological sequences). As a result, we presented an $O(n^3)$ time and $O(n^2)$ space algorithm to compute this mutation distance. Note that the mutation operations we considered allow both inversion and transposition and the lengths of two substrings exchanged by a transposition are not necessarily identical. It is worth mentioning that the algorithm we proposed in this study can be used to solve the approximate string matching problem under

non-overlapping inversion and/or transposition distance.

References

- A. Amir, T. Hartman, O. Kapah, A. Levy, E. Porat, On the cost of interchange rearrangement in strings, SIAM Journal On Computing 39 (2009) 1444–1461
- [2] D. Cantone, S. Cristofaro, S. Faro, Efficient string-matching allowing for non-overlapping inversions, Theoretical Computer Science 483 (2013) 85–95
- [3] D. Cantone, S. Faro, E. Giaquinta, Approximate string matching allowing for inversions and translocations, Proceedings of the Prague Stringology Conference 2010, pp. 37–51
- [4] D.-J. Cho, Y.-S. Han, H. Kim, Alignment with non-overlapping inversions on two strings, in: S.P. Pal, K. Sadakane (eds.) WALCOM 2014, LNCS, vol. 8344, Springer, Heidelberg, 2014, pp. 261–272.
- [5] D.-J. Cho, Y.-S. Han, H. Kim, Alignment with non-overlapping inversions and translocations on two strings, Theoretical Computer Science, in press (2014)
- [6] S. Grabowski, S. Faro, E. Giaquinta, String matching with inversions and translocations in linear average time (most of the time), Information Processing Letters 111 (2011) 516–520
- [7] S.-Y. Huang, C.-H. Yang, T.T. Ta, C.L. Lu, Approximate string matching problem under non-overlapping inversion distance, in: Workshop on Algorithms and Computation Theory in conjunction with 2014 International Computer Symposium, 2014, pp. 38–46
- [8] Y.-L. Huang, C.L. Lu, Sorting by reversals, generalized transpositions, and translocations using permutation groups. Journal of Computational Biology 17 (2010) 685–705
- [9] O. Kapah, G.M. Landau, A. Levy, N. Oz: Interchange rearrangement: the element-cost model, Theoretical Computer Science, 410 (2009) 4315–4326
- [10] F.T. Zohora, M.S. Rahman, Application of consensus string matching in the diagnosis of allelic heterogeneity, in: Basu, M., Pan, Y., Wang, J. (eds.) ISBRA 2014. LNBI 8492, vol. 8344, Springer, Switzerland, 2014, pp. 163–175