

# Computing Translocation Distance by a Genetic Algorithm

Lucas A. da Silveira<sup>†</sup>, José L. Soncco-Álvarez<sup>†</sup>, Thaynara A. de Lima<sup>‡</sup> and Mauricio Ayala-Rincón<sup>†</sup>

<sup>†</sup>Departament of Computer Science  
Universidade de Brasília  
70900-010 Brasília D.F., Brazil

<sup>‡</sup>Department of Mathematics  
Universidade Federal de Goiás — Campus II  
74690-900 Goiânia, Brazil

**Abstract**—Translocation is a useful operation on strings with challenging questions in combinatorics of permutations and interesting applications in analysis of sequences. A translocation operation essentially is the interchange of prefixes and suffixes among two substrings of a string. For the case of genomes represented as strings, symbols represent genes and chromosomes are modeled as substrings of the genomes; thus, translocation is an operation that models the interaction between chromosomes among a genome. The translocation distance between two genomes is defined as the minimum number of translocations to convert one genome into the other and had been proved to be a meaningful manner of modeling the evolutive distance between organisms. The particular case of unsigned genomes, those in which the orientation of the genes are not considered, is particularly difficult, while the signed case, in which the orientation of genes is considered, has been proved to be polynomially decidable. This paper compiles a proof of the  $\mathcal{NP}$ -hardness and presents an innovative GA approach to solve the unsigned translocation distance problem. As distinguished feature, the proposed GA uses as fitness function the translocation distance for randomly generated signed versions of the unsigned genomes. Experiments over randomly generated strings (synthetic chromosomes) show that the proposed GA approach compute answers that are better than those computed by an  $1.5+\varepsilon$ -approximation algorithm, the latter also implemented as part of this work.

## I. INTRODUCTION

The comparison of biological sequences is a problem of great relevance in the field of bioinformatics. By doing this comparison we can determine the evolutionary relationships between organism through the reconstruction of the sequence of evolutionary events that transform a genome into another. These rearrangements mechanisms include operations such as: reversals, transpositions, and translocations. The reversal and transposition operations are generally applied to genomes of only one chromosome ([1], [2], [3]), however translocations are operations that are applied over multiple chromosomes ([4], [5]).

### A. Related work

The unsigned translocation distance problem (UTD, for short) was proved  $\mathcal{NP}$ -hard by Zhu and Wang in [6] using previous results that relate the complexity of other problems such as decomposition in Eulerian cycles and color alternate cycle decomposition [7] and decomposition in  $k$ -Cliques [8]. As background this work compiles a complete proof.

The signed translocation distance problem results much more simpler than the unsigned case. Indeed, the orientation of genes provides a strong constraint in the genomes that

reduces drastically the combinatorics of the problem. The first polynomial  $O(n^3)$  solution was proposed by Hannenhalli in [4] giving rise to several other polynomial algorithms such as a quadratic one proposed in [9] and a linear algorithm proposed by Bergeron et al in [5] that has been implemented as part of the system UniMoG [10] and in the current work reimplemented in C in order to compute the fitness function of the proposed GA.

For unsigned genomes, that are the ones treated in this paper, known approximate solutions include the following. In [11], Kececioğlu and Ravi gave a ratio-2 approximation algorithm for computing the translocation distance between unsigned genomes; Cui et al. presented an  $1.75$ -approximation algorithm in [12], and further improved the approximation ratio to  $1.5+\varepsilon$  in [13]. Currently, to the best of our knowledge, the best approximation algorithm is one of ratio  $1.408+\varepsilon$  proposed in [14].

### B. Contribution

In this work we present a compilation of the proof of  $\mathcal{NP}$ -hardness of the unsigned translocation distance problem and a GA approach for solving this problem. The fitness function is based on linear computations of translocation distance for signed versions of the genomes as implemented in [5]. To verify the quality of the solutions computed by the GA, the  $1.5+\varepsilon$  approximation algorithm [13] was implemented as part of this work.

In the literature, the best proposed algorithm for UTD is the  $1.408+\varepsilon$ -approximation algorithm, however we use the algorithm  $1.5+\varepsilon$ -approximation, because the  $1.408+\varepsilon$ -approximation algorithm requires computing approximate solutions of the maximum set packing problem with set size at most 3 (which is  $\mathcal{NP}$ -complete [15]) and for the maximum independent set problem with maximum degree 4 (which is  $\mathcal{NP}$ -Complete [15]), and the computation of these problems can not be done in a straightforward manner. Moreover, for our requirements, the quality of the solutions provided by both approximation algorithms are similar since the ratio values are very closed. Thus, we implemented the  $1.5+\varepsilon$ -approximation algorithm and use it as a good mechanism to control the quality of the solutions of the proposed GA approach. This algorithm takes as inputs two unsigned genomes  $A$  and  $B$  (identity genome) and provides as output a signed genome  $\bar{A}$ . The idea of the algorithm is the following: compute the decomposition in cycles of the breakpoint graph  $G_u(A, B)$  (this will be defined in Section II), and from this decomposition attribute signals for the genes in  $A$  obtaining a signed genome

$\vec{A}$ ; then, one can calculate the translocation distance for  $\vec{A}$  and the identity genome using the linear time algorithm [5] for the UTD problem.

Several experiments were performed for calculating the translocations distance, indeed, sets of hundred genomes were randomly generated as input, with each set of genomes having lengths 10, 20, 30 until 150 genes. The results of these experiments showed that the proposed standard GA outperforms the quality of solutions computed by the  $1.5+\varepsilon$  approximation algorithm. Regarding running time, the GA takes only 10 seconds for genomes of length 150. The code was implemented in C and is available at [www.mat.unb.br/~ayala/publications.html](http://www.mat.unb.br/~ayala/publications.html).

### C. Organization

Initially, Section II presents the necessary background to understand the UTD and Section III presents a detailed compilation of the proof of  $\mathcal{NP}$ -hardness of the UTD. Afterwards, Section IV explains the  $1.5+\varepsilon$ -approximation algorithm and Section V introduces the standard GA for solving the UTD. Finally, before concluding and presenting future work in Section VIII, Sections VI and VII respectively present the experiments and discusses quality of the results.

## II. BACKGROUND

Standard definitions and notations are used (e.g. [5], [4], [13]).

### A. Genes, Chromosomes and Genomes

In order to represent the genes inside genomes of organisms, each gene is associated with an integer number. A signed integer represents an oriented gene and an unsigned integer a non oriented gene. A chromosome is a sequence of genes and a genome is constituted by a set of chromosomes. To simplify the model, we consider that each gene appears only once in the genome. So, a genome  $G$  with  $n$  oriented genes and  $N$  chromosomes can be seen as:

$$G = \{(x_{11} \dots x_{1r_1}), \dots, (x_{k1} \dots x_{kr_k}), \dots, (x_{N1} \dots x_{Nr_N})\}$$

where  $\sum_{i=1}^N r_i = n$ ,  $x_{ij} \in \{\pm 1, \dots, \pm n\}$  and  $|x_{ij}| \neq |x_{lk}|$  whenever  $i \neq l$  or  $j \neq k$ . For the unsigned version,  $x_{ij} \in \{1, \dots, n\}$ .

Chromosomes do not have orientation. Thus, the chromosomes  $X = (x_1, x_2, \dots, x_k)$  and  $X' = (-x_k, -x_{k-1}, \dots, -x_1)$  are the same in the signed case; whereas  $X$  and  $X'' = (x_k, x_{k-1}, \dots, x_1)$  are equal in the unsigned case. So, for example  $G = \{(+1 -3), (-4 +2 -5)\}$  is a genome with 5 genes and 2 chromosomes; furthermore,  $G$  and  $G' = \{(+1 -3), (+5 -2 +4)\}$  are the same genome.

For differentiate signed from unsigned genomes the former are denoted with an arrow:  $\vec{G}$ .

Genomes with a sole chromosome can be seen as permutations  $\pi$  in the symmetric group  $S_n$ . Indeed, a permutation is a bijective function from  $\{1, \dots, n\}$  into the same set of naturals. A permutation  $\pi$  can be represented as  $(\pi_1, \dots, \pi_n)$ , where  $\pi_i$  abbreviates  $\pi(i)$ , for  $1 \leq i \leq n$ .

### B. Sub-permutations

Let  $A$  and  $B$  be genomes with the same genes and  $S = (x_1, x_2, \dots, x_n)$  be a chromosome in  $A$ . A sub-permutation in the chromosome  $S$  in  $A$  to  $B$  (for short,  $SP$  in  $A$  to  $B$ ) is a segment  $[x_i, x_{i+1}, \dots, x_{i+l}]$  occurring in  $S$  with at least 3 genes, such that exactly the naturals between  $|x_i|$  and  $|x_{i+l}|$  occur in the set  $\{|x_{i+1}|, \dots, |x_{i+l-1}|\}$  and there is another segment  $[y_j, y_{j+1}, \dots, y_{j+l}]$  in some chromosome  $T$  of the genome  $B$ , satisfying:

- $|x_i| = |y_j|$  and  $|x_{i+l}| = |y_{j+l}|$ ;
- $\{|x_{i+1}|, \dots, |x_{i+l-1}|\} = \{|y_{j+1}|, \dots, |y_{j+l-1}|\}$ ;
- $[x_i, x_{i+1}, \dots, x_{i+l}] \neq [y_j, y_{j+1}, \dots, y_{j+l}]$ .

A *MinSP* is a  $SP$  in  $A$  to  $B$  that does not contain any other  $SP$ . For instance, consider the genomes:

$$A = \{(1, 3, 2, 4, 5, 8, 6), (7, 9)\} \text{ and}$$

$$B = \{(1, 2, 3, 4, 5, 6), (7, 8, 9)\};$$

$[1, 3, 2, 4, 5]$  is a  $SP$  and  $[1, 3, 2, 4]$  is a *MinSP*. For the signed genomes below,  $[+1, -3, +2, +4, +5]$  is a  $SP$  and  $[+1, -3, +2, +4]$  is a *MinSP*:

$$\vec{A} = \{(+1, -3, +2, +4, +5, +8, +6), (+7, +9)\} \text{ and}$$

$$\vec{B} = \{(+1, +2, +3, +4, +5, +6), (+7, +8, +9)\}.$$

### C. Breakpoint Graphs

Breakpoint graphs are an important data structure used in combinatorics of permutations and also useful for sorting genomes by translocations and other biological mutations.

Given a chromosome  $X = (x_1, x_2, \dots, x_n)$  of a (signed or unsigned) genome, we say that the genes  $x_i$  and  $x_{i+1}$ , for  $1 \leq i \leq n-1$ , are *adjacent*; otherwise, they are *not adjacent*. Also, genes in different chromosomes are not adjacent.

Consider two signed genomes  $\vec{A}$  and  $\vec{B}$  with the same genes and number of chromosomes. We can build the breakpoint graph  $G_s(\vec{A}, \vec{B})$  as follows. For all chromosomes  $X = (x_1, x_2, \dots, x_n)$  of  $\vec{A}$  and  $Y = (y_1, y_2, \dots, y_m)$  of  $\vec{B}$  the following elements are included:

- Vertices: a left-right ordered pair of vertices  $(l(x_i), r(x_i)) = (-x_i, +x_i)$ , for each gene  $x_i$ ,  $1 \leq i \leq n$ ;
- Edges: there is a black edge between  $r(x_i)$  and  $l(x_{i+1})$ , for  $1 \leq i < n$  and, there is a gray edge between  $+y_j$  and  $-y_{j+1}$ , if  $y_j$  and  $y_{j+1}$  are adjacent in  $B$ ,  $1 \leq j \leq m$ .

Figure 1 illustrates this notion.

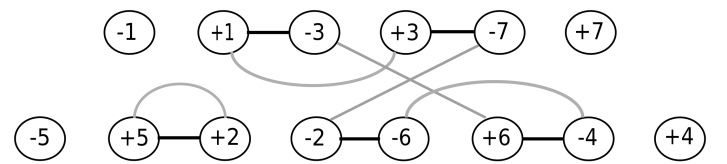


Fig. 1.  $G_s(\vec{A}, \vec{B})$  for  $\vec{A} = \{(+1, +3, +7), (+5, -2, +6, +4)\}$  and  $\vec{B} = \{(+1, -3, -6, +4), (+5, -2, +7)\}$

Notice one gray and one black edge are incident to each vertex in  $G_s(\vec{A}, \vec{B})$ , except for those vertices at the ends

of chromosomes. Therefore, breakpoint graphs can only be decomposed into color alternating cycles and univocally (cf. [7]). A cycle is called *long*, if it contains at least two black (or gray) edges, otherwise it is *short*.

Breakpoint graphs are also defined for the unsigned case. Consider two unsigned genomes  $A$  and  $B$  with the same genes and number of chromosomes. The breakpoint graph  $G_u(A, B)$  for  $A$  and  $B$  is constructed as follows: vertices are given by the genes in  $A$  and for all chromosomes  $X = (x_1, x_2, \dots, x_n)$  in  $A$ , there is a black edge between  $x_i$  and  $x_{i+1}$ ,  $1 \leq i < n$ ; and, for all chromosomes  $Y = (y_1, y_2, \dots, y_m)$  of  $B$  there is a gray edge between  $y_j$  and  $y_{j+1}$ , whenever  $y_j$  and  $y_{j+1}$  are adjacent in  $B$ . Figure 2 illustrates this notion.

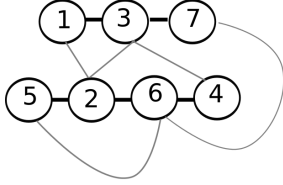


Fig. 2.  $G_u(A, B)$  for  $A = \{(1, 3, 7), (5, 2, 6, 4)\}$  and  $B = \{(1, 2, 3, 4), (5, 6, 7)\}$

The graph  $G_u(A, B)$  can be partitioned into a set of alternating cycles. Notice that, each vertex in  $G_u(A, B)$  has the same number of black and gray incident edges: vertices associated with genes at the end of chromosomes in  $A$  have only one black and one gray edge and internal genes have exactly two black and two gray incident edges. Thus, there is more than one way to partition  $G_u(A, B)$  into alternating cycles.

Breakpoint graphs will also be defined for permutations and we will see that these are almost those graphs obtained for genomes  $A$  and  $B$  as before, where  $A$  has only one chromosome.

#### D. Translocation

A translocation is said to be *active* in two chromosomes  $X$  and  $Y$  when both are cut and represented as  $X = (X_1, X_2)$  and  $Y = (Y_1, Y_2)$  and the segments produced on both chromosomes are interchanged, transforming  $X$  and  $Y$  in two new chromosomes  $X'$  and  $Y'$ . A translocation operation works with the assumption that the segments  $X_1, X_2, Y_1$  and  $Y_2$  are not empty.

In the translocation scenario, the literature presents two types of operations over segments of two chromosomes: *Prefix-Prefix* and *Prefix-Suffix*. Given two signed chromosomes  $X = (x_1, x_2, \dots, x_n)$  and  $Y = (y_1, y_2, \dots, y_m)$  in a genome, applying the translocation by *Prefix-Prefix*  $\rho(X, Y, x_i, y_k)$ , one obtains two new chromosomes  $X' = (x_1, \dots, x_i, y_{k+1}, \dots, y_m)$  and  $Y' = (y_1, \dots, y_k, x_{i+1}, \dots, x_n)$ . On the other hand, the translocation by *Prefix-Suffix*  $\theta(X, Y, x_i, y_k)$  produces the new chromosomes  $X' = (x_1, \dots, x_i, -y_k, \dots, -y_1)$  and  $Y' = (-y_m, \dots, -y_{k+1}, x_{i+1}, \dots, x_n)$  (See Figure 3).

**Example:** Consider the genome  $\vec{A} = \{X, Y, Z\}$  with  $X = (+1, +2, -7, +5)$ ,  $Y = (+4, +3)$  and  $Z = (+6, -8, +9)$ . The

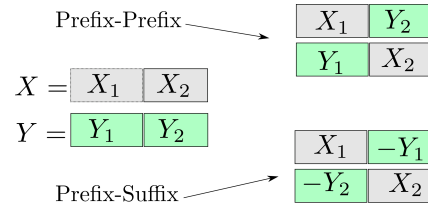


Fig. 3. *Prefix-Prefix* and *Prefix-Suffix* type translocations

translocation  $\rho(X, Y, +2, +4)$  transforms  $\vec{A}$  into the genome

$$\vec{A}' = \{(+1, +2, +3), (+4, -7, +5), Z\}.$$

Applying the translocation  $\theta(X, Z, +2, +6)$  to  $\vec{A}'$ , one obtains the genome

$$\vec{A}'' = \{(+1, +2, -6), Y, (-9, +8, -7, +5)\}.$$

For the unsigned case, translocation is defined as follows. Consider two unsigned chromosomes  $X = (x_1, x_2, \dots, x_n)$  and  $Y = (y_1, y_2, \dots, y_m)$ . A translocation by *Prefix-Prefix*  $\rho(X, Y, x_i, y_k)$  transforms  $X$  and  $Y$  into two new chromosomes  $X' = (x_1, \dots, x_i, y_{k+1}, \dots, y_m)$  and  $Y' = (y_1, \dots, y_k, x_{i+1}, \dots, x_n)$ ; whereas a translocation by *Prefix-Suffix*  $\theta(X, Y, x_i, y_k)$  transforms  $X$  and  $Y$  into  $X' = (x_1, \dots, x_i, y_k, \dots, y_1)$  and  $Y' = (y_m, \dots, y_{k+1}, x_{i+1}, \dots, x_n)$ .

**Example:** Consider the unsigned genome  $A = \{X, Y, Z\}$ , with chromosomes  $X = (1, 2, 7, 5)$ ,  $Y = (4, 3)$  and  $Z = (6, 8, 9)$ .  $\rho = (X, Y, 2, 4)$  transforms  $A$  into

$$A' = \{(1, 2, 3), (4, 7, 5), Z\}.$$

$\theta(X, Z, 2, 6)$  transform  $A$  into

$$A'' = \{(1, 2, 6), Y, (9, 8, 7, 5)\}.$$

Given a chromosome  $X = (x_1, x_2, \dots, x_k)$ , the genes  $x_1$  and  $-x_k$  are called *tails* of  $X$ . Two genomes are called *co-tails* if their sets of *tails* are the same. The genomes  $\vec{B}$  and  $\vec{C}$  below are co-tails since they share the same set of tails, that is  $\{+1, -6, +7, -10\}$ .

$$\vec{B} = \{(+1, +2, -4, +3, +5, +6), (+7, -9, +8, +10)\},$$

$$\vec{C} = \{(+1, +2, +3, +4, +5, +6), (+7, +8, +9, +10)\}.$$

Notice also that genomes  $\vec{A}, \vec{A}'$  and  $\vec{A}''$  as well as  $A, A'$  and  $A''$  of previous example are co-tails.

This property is essential, when we consider genome rearrangement through translocations, because translocations by *Prefix-Prefix* and *Prefix-Suffix* do not alter the set of *tails* of a genome. So, in order to transform the genome  $A$  into  $B$  by translocations the following conditions must be satisfied: the number of chromosomes and genes of  $A$  and  $B$  must be the same and  $A$  and  $B$  must be co-tails.

In the rest of the paper, unless otherwise stated, we will consider only unsigned genomes.

### E. Translocation distance

We are interested in studying the problem of sorting genomes by translocations. The problem can be described as follows: consider two unsigned genomes  $A$  and  $B$  with  $n$  genes, where the genes of the genome  $B$  are in increasing order and  $A$  and  $B$  are co-tails. Our goal is to find a sequence  $\delta_1, \delta_2, \dots, \delta_t$  of translocations that transform  $A$  into  $B$ , and  $t$  is minimum; the number  $t$  is called the translocation distance between  $A$  and  $B$ . For the signed case, the problem is defined analogously; but, the genes of the genome  $B$  are positive, furthermore they are in increasing order.

The complexity of the translocation distance problem is related with the maximum decomposition into alternating cycles of breakpoint graphs. Since there is only one possible decomposition into alternating cycles of the breakpoint graph of signed genomes, the translocation distance problem results of polynomial complexity; whereas for the unsigned version, the problem is  $\mathcal{NP}$ -hard as will be explained in the next section.

### III. UNSIGNED TRANSLOCATION DISTANCE IS NP-HARD

The original proof of this fact is in [6]; here, we will compile all necessary steps providing a self-contained proof.

#### A. Screenplay of the Proof

For a good understanding of the stages of the proof, a screenplay will be presented containing all necessary steps that were implicitly or explicitly used in [6] and the associated references. Some problems should be defined first.

**$k$ -cliques:** check if the set of edges of a graph  $H$  can be partitioned into cliques of size  $k$ . For  $k \geq 3$ ,  $k$ -cliques is known to be an  $\mathcal{NP}$ -complete problem.

**MAX-ECD:** consider an Eulerian graph  $H$ . The problem consists in finding a maximum decomposition in cycles of  $H$ , i.e., a partition of the set of edges of  $H$  in the maximum number of cycles.

**MAX-ACD:** given a breakpoint graph  $G(\pi)$  of a permutation  $\pi$ , the problem consists in finding a maximum decomposition in alternating cycles of  $G(\pi)$ .

The proof is organized in the following steps (See Figure 4):

- Initially, Section III-B demonstrates that the problem of graph partitioning in cycles for  $k = 3$  is an instance of the  $MAX-ECD$  problem [8]; thus,  $MAX-ECD$  is an  $\mathcal{NP}$ -complete problem.
- Afterwards, Section III-C presents a polynomial reduction from  $MAX-ECD$  into  $MAX-ACD$  [7]; consequently,  $MAX-ACD$  is a  $\mathcal{NP}$ -hard problem;
- Finally, Section III-D polynomially reduces  $MAX-ACD$  into the translocation distance problem [6]; thus, the translocation distance problem is  $\mathcal{NP}$ -hard.

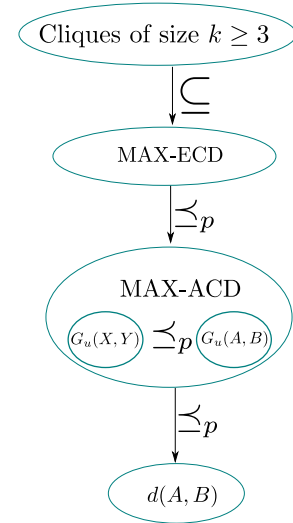


Fig. 4. Reductions for  $\mathcal{NP}$ -hardness of unsigned translocation distance

#### B. $k$ -cliques $\subseteq$ $MAX-ECD$

In the early 1980's Holyer proved that partitioning a graph in  $k$ -cliques of size  $k$ , for  $k \geq 3$  is an  $\mathcal{NP}$ -complete problem [8]. In particular, for  $k = 3$  one wants to check if the edge set of the graph can be partitioned into triangles. In this case, the graph can be assumed Eulerian. Thus, the problem of finding a partition of the set of edges of an Eulerian graph into triangles is  $\mathcal{NP}$ -complete. Furthermore, since the decomposition of a graph into triangles gives the maximum Eulerian decomposition, one can conclude that 3-cliques  $\subseteq$   $MAX-ECD$ .

#### C. $MAX-ECD \preceq_p$ $MAX-ACD$

In this section we present details of the polynomial reduction from  $MAX-ECD$  to  $MAX-ACD$  proposed by Caprara in [7]. Before this, some definitions and properties must be given.

1) **Breakpoint Graph  $G(\pi)$ :** Consider a permutation  $\pi = (\pi_1, \dots, \pi_n)$  in  $S_n$  representing a genome constituted by only one chromosome.

The breakpoint graph  $G(\pi) = \langle V, E = R \cup B \rangle$  of  $\pi$  is constructed as follows: initially, add to  $\pi$  the elements  $\pi_0 = 0$  and  $\pi_{n+1} = n + 1$ , and consider  $\pi' = (0, 1, \dots, n + 1)$ . Each node  $v \in V$  represents an element of  $\pi'$ .

The breakpoint graph  $G(\pi)$  is a bi-colored graph, where the set of edges  $E$  is partitioned into two subsets: red and blue edges  $R$  and  $B$ . There is a red edge  $(\pi_i, \pi_{i+1})$  whenever  $|\pi_i - \pi_{i+1}| \neq 1$ , for  $0 \leq i \leq n$ , i.e., there is a red edge between consecutive vertices  $\pi_i$  and  $\pi_{i+1}$  that have non consecutive labels. In this case, the pair  $\pi_i, \pi_{i+1}$  is called a *breakpoint* of  $G(\pi)$ . There is a blue edge between vertices with labelled with  $i$  and  $i + 1$ ,  $1 \leq i \leq n$ , whenever  $i$  and  $i + 1$  are not in consecutive positions in  $\pi$ . The Figure 5 illustrates the graph  $G(\pi)$  for the genome  $\pi = (2, 4, 1, 3)$ .

**Theorem 1** (Properties of  $G(\pi)$  - Th. 4 in [7]). *A bi-colored graph  $G = \langle V, R \cup B \rangle$  is the breakpoint graph of some genome  $\pi$  iff*

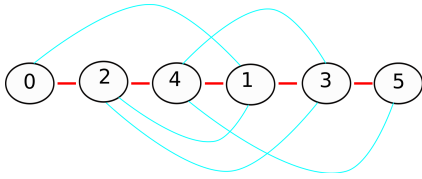


Fig. 5. Breakpoint graph  $G(\pi)$  associated to  $\pi = (2, 4, 1, 3)$

- Each connected component of the subgraphs  $G(R)$  and  $G(B)$ , induced by the red and blue edges resp., is a simple path;
- Each node  $i \in V$  has the same degree (0, 1 or 2) in  $G(R)$  and  $G(B)$ ;
- There are no edges in  $G(R)$  and  $G(B)$  connecting the same vertices.

Sufficiency follows from definition of breakpoint graphs. Necessity requires construction of a permutation  $\pi$  through Hamiltonian matchings [7].

An alternating cycle of  $G(\pi)$ , is a sequence of edges  $r_1, b_1, r_2, b_2, \dots, r_m, b_m$ , where  $r_i \in R$ ,  $b_i \in B$  for  $i = 1, \dots, m$ ;  $r_i$  and  $b_j$  have a common node for  $i = j = 1, \dots, m$  and for  $i = j + 1$ ,  $j = 1, \dots, m$  (where  $r_{m+1} = r_1$ ); and  $r_i \neq r_j$ ,  $b_i \neq b_j$  for  $1 \leq i < j \leq m$ .

A decomposition in alternating cycles of  $G(\pi)$  is a collection of alternating disjoint edge cycles, such that each edge of  $G(\pi)$  is contained in exactly one cycle of the collection. Thus, *MAX-ACD* is an optimization problem that consists in searching a maximum decomposition of  $G(\pi)$  in alternating cycles. For instance see the *MAX-ACD* in the Figure 6.

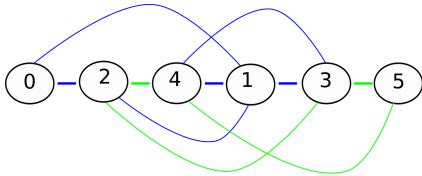


Fig. 6. *MAX-ACD* of size 2 for  $G(\pi)$  for  $\pi = (2, 4, 1, 3)$  (Fig. 5)

2) *MAX-ACD* is  $\mathcal{NP}$ -Hard: The proof is based on a polynomial reduction from *MAX-ECD* to *MAX-ACD*. Given an Eulerian graph  $H = \langle W, E \rangle$  one can build a bi-colored graph  $G$  in polynomial time, such that there exists a one-to-one correspondence between cycles of  $H$  and alternating cycles of  $G$ . The graph  $G$  is built replacing each node  $v$  of  $H$  of degree  $d$ , by a bi-colored subgraph  $G(v)$ , containing  $d$  bottom nodes among other vertices, where  $d$  is the degree of  $v$  (see Figure 7). All connections between the original vertices of  $H$  remain represented by red connections between different base nodes of the bi-colored subgraphs  $G(v) \in G$ . To complete the proof, the graph  $G$  must satisfy the characterization of Theorem 1.

The internal structure of each subgraph  $G(v) \in G$  has  $d$  base nodes and  $m$  levels, abbreviated as  $G(d, m)$ . Let  $s := \frac{d}{2}$  and  $r := \lceil \frac{d}{4} \rceil$ .

Each level  $l$ ,  $l = 1, \dots, m$  contains  $2s + 1$  nodes;  $s + 1$  of them are top-level nodes, denoted by  $q_1^l, \dots, q_{s+1}^l$ , and

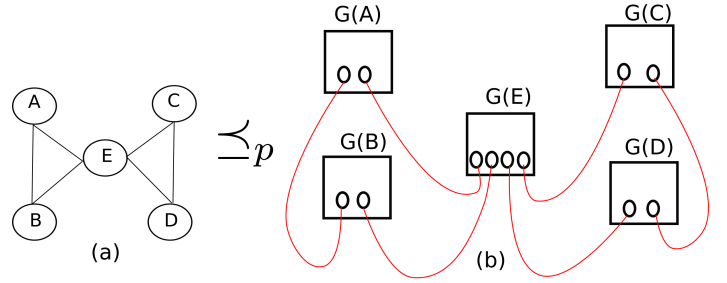


Fig. 7. (a) The Eulerian graph  $H$ . (b) The bi-colored graph  $G(\pi)$  constructed from  $H$ .

the other  $s$  are lower-level nodes, denoted by  $p_1^l, \dots, p_s^l$ . Top-level nodes  $q_1^l, \dots, q_{s+1}^l$  are connected to lower-level nodes  $p_1^l, \dots, p_s^l$  by  $d$  red edges  $(q_i^l, p_i^l), (q_{i+1}^l, p_i^l)$ , for  $i = 1, \dots, s$ .

Also, for  $l = 1, \dots, m - 1$  the top-level nodes  $q_1^l, \dots, q_{s+1}^l$  are connected to the lower-level nodes of the level  $l + 1$  by  $d$  blue edges  $(p_i^{l+1}, q_i^l), (p_i^{l+1}, q_{i+1}^l)$ , for  $i = 1, \dots, s$ .

Finally, the upper nodes of the last level  $m$  are connected to each other by  $s$  blue edges  $(q_i^m, q_{i+1}^m)$ , for  $i = 1, \dots, s$ , and the  $d$  bottom nodes, denoted by  $b_1, \dots, b_d$ , are connected to the lower-level nodes of the first level by  $d$  blue edges  $(b_{2i-1}, p_i^1), (b_{2i}, p_i^1)$ , for  $i = 1, \dots, s$  (see Figure 8).

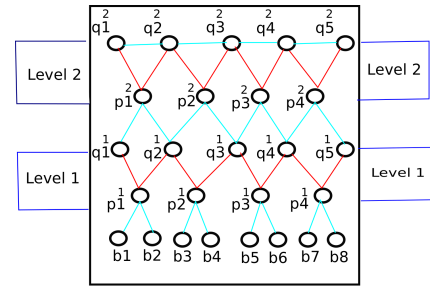


Fig. 8. Subgraph  $G(d, m)$  for  $d = 8$  and  $m = 2$ .

Observe that all nodes of  $G(d, m)$  except bottom nodes have the same number of red and blue incident edges. Given two nodes of  $G(d, m)$  an alternating path between such nodes is a path where the colors of the edges are alternate. Given any two different bottom nodes  $b_i$  and  $b_j$ , it is always possible to build an alternating path between  $b_i$  and  $b_j$ , whenever we have enough number of levels in  $G(d, m)$ . Indeed, by Theorem 5 in [7], the edge set of  $G(d, m)$  can be decomposed into  $s$  alternating paths, connecting any selection of different pairs of bottom nodes, iff  $m \geq r(s - 1) + 1$ . Thus, one can guarantee that it is possible from a bottom node  $b_i$  achieve a different bottom node  $b_j$  by an alternating path and then connect  $b_j$  to another subgraph belonging to  $G$  by a red edge (see Figure 7).

Thus, if there is a cycle in  $H$  then it can be represented by an alternating cycle in  $G$  and vice-versa; consequently there is a correspondence between a cycle decomposition of  $H$  and  $G$ . The Figure 9 illustrates this correspondence for a particular case.

To conclude, one can observe that  $G$  satisfies the conditions of the Theorem 1. Consequently,  $G$  is a breakpoint graph. The reduction is done in polynomial time choosing  $m = r(s - 1) +$

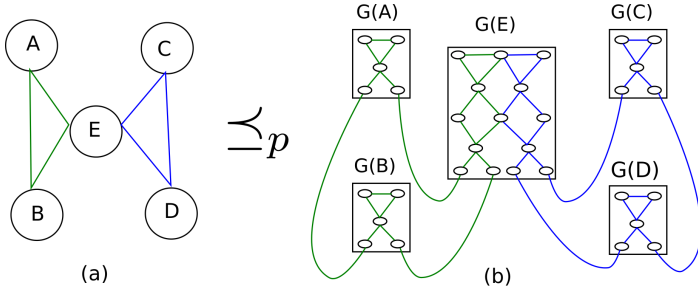


Fig. 9. (a) Eulerian graph H containing two cycles. (b) The bi-colored graph G, representing the same two cycles of H.

1. Thus, we have a polynomial time reduction from *MAX-ECD* to *MAX-ACD*, and *MAX-ACD* is  $\mathcal{NP}$ -hard.

#### D. *MAX-ACD* $\leq_p$ *UTD*

In this section, a polynomial reduction from *MAX-ACD* to *UTD* is presented following the presentation in [6]. This allows one concluding that the latter problem is  $\mathcal{NP}$ -hard.

Let  $X$  and  $Y$  be two unsigned chromosomes. Without loss of generality, let  $X = (g_1, g_2, \dots, g_n)$  and  $(Y = 1, 2, \dots, n)$ , where  $\{g_1, g_2, \dots, g_n\} = \{1, 2, \dots, n\}$  and  $g_1 = 1, g_n = n$ . From  $X$  and  $Y$ , one can build two genomes  $A = \{X_1, X_2\}$  and  $B = \{Y_1, Y_2\}$  as will be described. Also, consider an integer  $d$  that will be used to control the amount of short cycles in the decomposition of  $G_u(A, B)$ ; this number will be detailed in Lemma 2. The chromosome  $X_1$  of the genome  $A$  is constructed by inserting  $n - 1$  new genes between two adjacent genes in  $X$  as follows:

$$X_1 = (1, t_{1,1}, g_2, t_{1,2}, \dots, g_{n-1}, t_{1,n-1}, n)$$

where,  $t_{1,k} = 3n - 2 + k, 1 \leq k \leq n - 1$ .

$X_2$  contains two types of new genes,  $t_{2,l} = n + l, 1 \leq l \leq 2(n - 1)$  and  $s_i = 4n - 3 + i, 1 \leq i \leq (n - 2)d$ .

$$\begin{aligned} X_2 = & (t_{2,1}, t_{2,2}, s_1, s_2, \dots, s_d, \\ & t_{2,3}, t_{2,4}, s_{d+1}, s_{d+2}, \dots, s_{2d}, \\ & \vdots \\ & t_{2,2(n-2)-1}, t_{2,2(n-2)}, s_{(n-3)d+1}, \dots, s_{(n-2)d}, \\ & t_{2,2(n-1)-1}, t_{2,2(n-1)}) \end{aligned}$$

To construct the genome  $B = \{Y_1, Y_2\}$ , consider the same integers  $t_{1,k}, t_{2,l}$  and  $s_i, 1 \leq k \leq n - 1, 1 \leq l \leq 2(n - 1), 1 \leq i \leq (n - 2)d$ , as calculated in  $A$ . The chromosome  $Y_1 = Y = (1, 2, \dots, n)$  and  $Y_2$  is built from  $X_2$  inserting  $t_{1,k}$  between  $t_{2,2k-1}$  and  $t_{2,2k}$  in  $X_2$ .

$$\begin{aligned} Y_2 = & (t_{2,1}, t_{1,1}, t_{2,2}, s_1, s_2, \dots, s_d, \\ & t_{2,3}, t_{1,2}, t_{2,4}, s_{d+1}, \dots, s_{2d}, \\ & \vdots \\ & t_{2,2(n-2)-1}, t_{1,n-2}, t_{2,2(n-2)}, s_{(n-3)d+1}, \dots, s_{(n-2)d}, \\ & t_{2,2(n-1)-1}, t_{1,n-1}, t_{2,2(n-1)}) \end{aligned}$$

At the end of the construction, each one of the genomes  $A$  and  $B$  has a total number of  $4n - 3 + (n - 2)d$  genes.

**Example.** Let  $X = (1, 3, 4, 2, 5)$  and  $Y = (1, 2, 3, 4, 5)$ ; the Figure 10 (a) illustrates the graph  $G_u(X, Y)$ . Consider  $d = 4$ . So, the genomes  $A$  and  $B$  are:

$$\begin{aligned} A = & \{X_1, X_2\}, \text{ where} \\ X_1 = & (1, 14, 3, 15, 4, 16, 2, 17, 5) \text{ and} \\ X_2 = & (6, 7, 18, 19, 20, 21, 8, 9, 22, 23, 24, \\ & 25, 10, 11, 26, 27, 28, 29, 12, 13) \end{aligned}$$

and

$$\begin{aligned} B = & \{Y_1, Y_2\}, \text{ where} \\ Y_1 = & (1, 2, 3, 4, 5) \text{ and} \\ Y_2 = & (6, 14, 7, 18, 19, 20, 21, 8, 15, 9, 22, 23, 24, \\ & 25, 10, 16, 11, 26, 27, 28, 29, 12, 17, 13) \end{aligned}$$

The graph  $G_u(A, B)$  is shown in the Figure 10 (b).

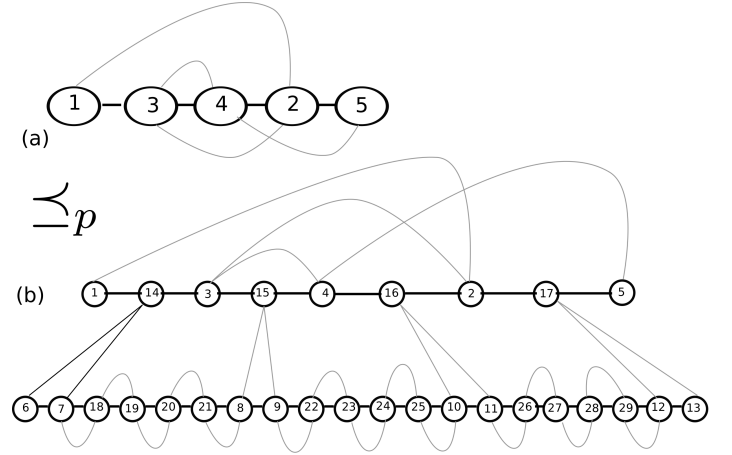


Fig. 10. Breakpoint graphs (a)  $G_u(X, Y)$  and (b)  $G_u(A, B)$

The Lemmas 2 and 5 clarify the relationship between a maximum decomposition into alternating cycles of  $G_u(X, Y)$  and the translocation distance between genomes  $A$  and  $B$ .

**Lemma 2** (Lemma 4 in [6]). *Assume  $d \geq n - 1$ . There is a decomposition of  $G_u(X, Y)$  into  $J$  alternating cycles iff there is a decomposition of  $G_u(A, B)$  into at least  $(n - 2)(d + 1) + J$  alternating cycles.*

*Proof:* A few details are aggregated to the original presentation given in [6].

**Sufficiency:** Assume that there is a decomposition  $M$  of  $G_u(X, Y)$  into  $J$  alternating cycles. For each cycle  $C \in M, C$  can be represented as a list  $C = u_1, u_2, \dots, u_{2k-1}, u_{2k}$ , where  $(u_{2i-1}, u_{2i})$  is a black edge and  $(u_{2i}, u_{2i+1})$  is a gray edge,  $1 \leq i \leq k$  and  $u_{2k+1} = u_1$ . Notice that if  $u_{2i-1} = g_j$  then  $u_{2i} = g_{j+1}$  or  $u_{2i} = g_{j-1}$ .

One can obtain a new cycle  $C'$  in  $G_u(A, B)$  replacing each black edge  $(u_{2i-1}, u_{2i})$  of  $C$  by the alternating path  $P_{2i-1, 2i}$ , where,

$$\begin{aligned} P_{2i-1, 2i} = & g_j, t_{1,j}, t_{2,2j-1}, t_{2,2j}, t_{1,j}, g_{j+1} \\ & \text{if } u_{2i-1} = g_j, u_{2i} = g_{j+1} \end{aligned}$$

$$\begin{aligned} P_{2i-1, 2i} = & g_j, t_{1,j-1}, t_{2,2j-3}, t_{2,2j-2}, t_{1,j-1}, g_{j-1} \\ & \text{if } u_{2i-1} = g_j, u_{2i} = g_{j-1}. \end{aligned}$$

So, to each cycle  $C \in M$  one can associate univocally a long cycle  $C'$  of  $G_u(A, B)$ ; such cycle is long because each alternating path  $P_{2i-2i}$  has 3 black and 2 gray edges. Notice that, the only edges of  $G_u(A, B)$  not used to build the  $J$  long cycles are the edges  $(t_{2,2i}, s_{(i-1)d+1}), \dots, (s_{id-1}, s_{id}), (s_{id}, t_{2,2i+1})$  for  $i \in \{1, \dots, n-2\}$ ; such edges can form  $(n-2)(d+1)$  short cycles. Consequently, the decomposition  $M$  of  $G_u(X, Y)$  into  $J$  alternating cycles induces a decomposition of  $G_u(A, B)$  into  $(n-2)(d+1) + J$  alternating cycles.

**Necessity:** Let  $M'$  be a set of  $(n-2)(d+1) + J$  alternating cycles forming a decomposition of  $G_u(A, B)$ . Notice that, only the edges  $(t_{2,2i}, s_{(i-1)d+1}), \dots, (s_{id-1}, s_{id}), (s_{id}, t_{2,2i+1})$  for  $i \in \{1, \dots, n-2\}$  can form short cycles of  $G_u(A, B)$ ; consequently, there are at most  $(n-2)(d+1)$  short cycles in  $M'$ . Also, an alternating cycle using a vertice in  $X_1$  must contain at least two black edges containing only vertices of  $X_1$  and at least two gray edges with a vertice in  $X_1$  and another in  $X_2$  (this implies that every cycle containing a vertice in  $X_1$  is a long cycle). Consequently, since  $X_1$  has  $2n-1$  vertices, there are at most  $n-1$  alternating cycles in  $M'$  using vertices in  $X_1$ . For any  $i$ , the  $d+2$  consecutive vertices  $(t_{2,2i}, s_{(i-1)d+1}), \dots, (s_{id-1}, s_{id}), (s_{id}, t_{2,2i+1})$  in  $X_2$  can not be in a same cycle in  $M'$ ; indeed, otherwise, the number of short cycles will be reduced by  $d+1 \geq n$ , by hipotesis, and the maximum number of cycles in  $M'$  would be

$$\begin{aligned} (n-3)(d+1) + (n-1) &= (n-2-1)(d+1) + n-1 \\ &\leq (n-2)(d+1) - n + n-1 \\ &< (n-2)(d+1) + J \end{aligned}$$

Thus, if a cycle in  $M'$  containing one of the gray edges  $(t_{1,j}, t_{2,2j-1})$  and  $(t_{1,j}, t_{2,2j})$ , it must contain both of them; otherwise, it would exist a cycle in  $M'$  using  $d+2$  consecutive vertices in  $X_2$ . Furthermore, if a long cycle in  $M'$  uses two consecutive vertices in the sequence  $(t_{2,2i}, s_{(i-1)d+1}), \dots, (s_{id-1}, s_{id}), (s_{id}, t_{2,2i+1})$ , such cycle must contain both of the black and gray edges between these two vertices; consequently, such long cycle can be decomposed in order to increase the number of short cycles. Thus, we can assume that all the  $(n-2)(d+1)$  short cycles are in  $M'$ .

Finally, because any long cycle from  $G_u(A, B)$  can not only contain gray edges in  $X_1$ , each long cycle in  $M'$  must take the path  $P_{1,2}, \dots, P_{2k-1,2k}$ , where  $P_{2i-1,2i} = u_{2i-1}, t_{1,j}, t_{2,2j-1}, t_{2,2j}, t_{1,j}, u_{2i}$ , and  $\{u_{2i-1}, u_{2i}\} = \{g_j, g_{j+1}\}$ . Replacing the path  $P_{2i-1,2i}$  by the black edge  $(u_{2i-1}, u_{2i})$ , one obtains an alternating cycle in  $G_u(X, Y)$ . Thus  $G_u(X, Y)$  can be decomposed into  $J$  alternating cycles. ■

Before Lemma 5, it is convenient to enunciate two theorems introduced in [6]. The first one relates translocation distance and cycles in the decomposition of breakpoint graphs and the second one translocation distance between unsigned genomes and their signed versions.

**Theorem 3** (Th. 1 in [6]). *The translocation distance between two signed genomes  $\vec{A}$  and  $\vec{B}$  satisfies  $d(\vec{A}, \vec{B}) \geq n - m - c_{AB}$ . If there is no sub-permutations for  $\vec{A}$  and  $\vec{B}$ , the  $d(\vec{A}, \vec{B}) = n - m - c_{AB}$ , where  $n$  is the number of genes,  $m$  the number of chromosomes and  $c_{AB}$  the number of cycles in the cycle decomposition of the breakpoint graph  $G_s(\vec{A}, \vec{B})$ .*

**Theorem 4** (Th. 2 in [6]). *Let  $A$  and  $B$  be unsigned genomes. Consider  $\vec{B}$  the signed genome obtained from  $B$  by setting every gene as positive. Then,  $d(A, B) = \min_{\vec{A} \in Spin(A)} d(\vec{A}, \vec{B})$ , where  $Spin(A)$  is the set of all signed genomes obtained from  $A$  by setting signs to its genes.*

**Lemma 5** (Lemma 5 in [6]). *There is a decomposition of  $G_u(A, B)$  into  $(n-2)(d+1) + J$  alternating cycles iff  $d(A, B) \leq 3n - 3 - J$ .*

*Proof: Sufficiency:* If  $G_u(A, B)$  can be decomposed into  $(n-2)(d+1) + J$  alternating cycles, there is a decomposition  $M'$  in at least  $(n-2)(d+1) + J$  alternating cycles of  $G_u(A, B)$ , where every long cycle has a gray edge containing a vertice of  $X_1$  and a vertice of  $X_2$ . From this decomposition it is possible construct signed genomes  $\vec{A}$  and  $\vec{B}$ , such that the number of alternating cycles of  $G_s(\vec{A}, \vec{B})$  is the same that in  $M'$ ; and because every long cycle in  $M'$  has a gray edge with a vertice in  $X_1$  and a vertice in  $X_2$ , there are no sub-permutations for  $\vec{A}$  and  $\vec{B}$ .

Let  $n$  the number of genes in both genomes  $A$  and  $B$ ,  $N$  the number of chromosomes and  $c_{AB}$  the number of alternating cycles in the decomposition  $M'$  of  $G_u(A, B)$ . So, by Theorem 4,  $d(A, B) \leq d(\vec{A}, \vec{B})$  and by Theorem 3,

$$\begin{aligned} d(\vec{A}, \vec{B}) &= n - N - c_{AB} \\ &= 4n - 3 + (n-2)d - 2 - \\ &= ((n-2)(d+1) + J) \\ &= 3n - 3 - J. \end{aligned} \tag{1}$$

Thus,  $d(A, B) \leq 3n - 3 - J$ .

**Necessity:** Let  $d(A, B) \leq 3n - 3 - J$ . By Theorem 4, there exists genomes  $\vec{A}$  and  $\vec{B}$  obtained from  $A$  and  $B$  such that  $d(\vec{A}, \vec{B}) \leq 3n - 3 - J$ .

Using again the Theorem 1,  $d(\vec{A}, \vec{B}) \geq 4n - 3 + (n-2)d - 2 - c_{\vec{A}\vec{B}}$ . Thus,

$$\begin{aligned} c_{\vec{A}\vec{B}} &\geq 4n - 3 + (n-2)d - 2 - d(\vec{A}, \vec{B}) \\ &\geq 4n - 3 + nd - 2d - 2 - 3n + 3 + J \\ &= n + nd - 2d - 2 + J \\ &= (n-2)(d+1) + J. \end{aligned} \tag{2}$$

Therefore, the maximum number of alternating cycles in a decomposition of  $G_u(A, B)$  is at least  $(n-2)(d+1) + J$ . ■

By Lemmas 2 and 5, there is a decomposition into  $J$  alternating cycles of  $G_u(X, Y)$ , if and only if, the translocation distance between  $A$  and  $B$  is at most  $3n - 3 - J$ .

Consider an instance of *MAX-ACD* containing  $n$  genes and  $d = n - 1$ . So, the corresponding instance of unsigned translocation distance has  $4n - 3 + (n-2)(n-1) = n^2 + n - 1$  genes.

Thus, there is a polynomial reduction from *MAX-ACD* to *Unsigned translocation distance problem* and the latter is  $\mathcal{NP}$ -hard.

#### IV. $1.5+\varepsilon$ -APPROXIMATE SOLUTION FOR UTD

In the search for approximate solutions for the translocation distance problem between unsigned genomes, Zhu and Wang noted that given an unsigned genome  $A$  (since  $B$  is considered as the identity genome, there is no need to give attention to it), depending on signals attributed to the genes, the minimum translocations necessary to order signed versions of  $A$  can vary. See a simple example in Figure 11. Thus, the solutions known in the literature to the unsigned case exploit this statement, applying complex heuristics, in order to acquire good approximate solutions.

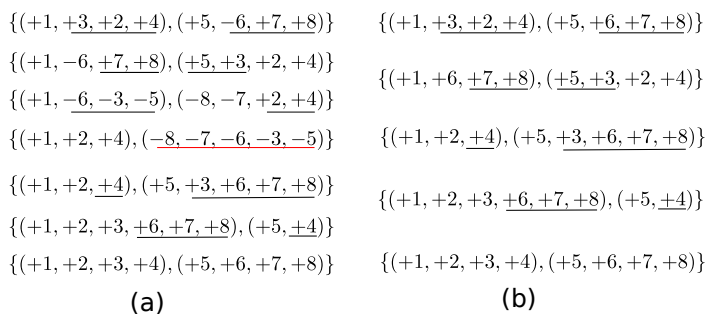


Fig. 11. The red line represents an inversion of chromosome and black lines translocations. (a) Genomes  $A = \{(+1, +3, +2, +4), (+5, -6, +7, +8)\}$  and  $B = \{(+1, +2, +3, +4), (+5, +6, +7, +8)\}$ . The translocation distance is 5. (b) Genomes  $A = \{(+1, +3, +2, +4), (+5, +6, +7, +8)\}$  and same  $B$ . The translocation distance is 4.

Here we described details of the implementation of the algorithm introduced in [13] that provides approximate solutions of ratio  $1.5+\varepsilon$  for translocation distance problem in the unsigned version. Solutions given by this implementation are used as a quality control for the solutions provided by the proposed GA. The strategy of this approximation algorithm consists in computing the decomposition in cycles of  $G_u(A, B)$ , and from this decomposition attribute signals for the genes in  $A$  obtaining some  $\bar{A}$ ; then, one can calculate the translocation distance for  $\bar{A}$  using a linear time algorithm for the signed version of the translocation distance problem.

##### A. Heuristics Used in the Approximation Algorithm

In the search for a decomposition in cycles of  $G_u(A, B)$ , all 1-cycles are maintained, where 1-cycles are formed by a black and a gray edge, such that appropriate signals are attributed to the genes involved in order to form an 1-cycle.

After obtaining the maximum number of 1-cycles, one seeks the maximum number of 2-cycles in polynomial time. A match graph  $F_{AB}$  of the breakpoint graph  $G_u(A, B)$  is constructed as follows:

**1** for each black edge in  $G_u(A, B)$  with at least one unsigned vertex, create a vertex in  $F_{AB}$ ;

**2** for each two vertices in  $F_{AB}$ , an edge is created connecting them if the two black edges in  $G_u(A, B)$  form a 2-cycle.

Let  $V$  and  $E$  be the set vertices and edges of  $F_{AB}$  respectively. A maximum match of  $F_{AB}$  is a set  $M \subseteq E$  such that:  $\forall v \in V, v$  has at most one edge incident in  $M$ .

Each edge in  $M$  represents a 2-cycle in  $G_u(A, B)$ . A 2-cycle in  $M$  is isolated if it does not share any edge with any

other 2-cycle. Otherwise, the 2-cycle is related. Since, a 2-cycle has two gray edges, it relates at most two 2-cycles.

A *related component*  $U$  consists of related 2-cycles  $c_1, c_2, \dots, c_k$ , where  $c_i$  is related with  $c_{i-1}$  ( $2 \leq i \leq k$ ), and each 2-cycle is not related to any other 2-cycle out of  $U$ . Also, a related component involves at most two chromosomes and can be only of one of the four types in Figure 12.

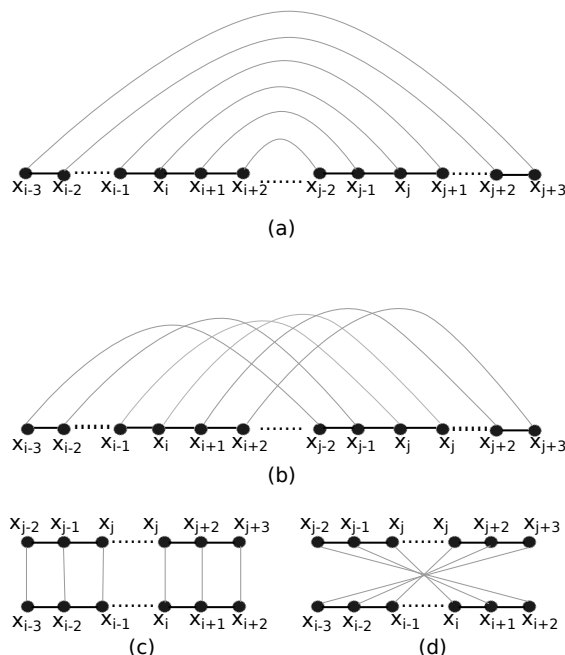


Fig. 12. The four types of related components containing 2-cycles. (a) e (b) The component is only at one chromosome. (c) e (d) Two chromosomes are involved in the component.

After the components are built, the focus will be the isolated components. From isolated 2-cycles it is possible to identify a special type of 2-cycle, called *simple minSP* (here, *minSP* coincides with the notion defined in Section II) or in the *short form SMSP*, those 2-cycles appear with their gray edges at the extremity twisted, and with all internal cycles with their black edges involved in 1-cycles. Let  $I_s = x_i, x_{i+1}, \dots, x_j$  an *SMSP*. If there exists a gray edge  $(x_{i-1}, x_j)$  or  $(x_i, x_{j+1})$  in  $G_u(A, B)$ , one can create a gray edge  $(r(x_{i-1}), r(x_j))$  or  $(l(x_i), l(x_{j+1}))$ , transforming  $I_s$  into a *removable SMSP (RSMSP)*. The *RSMSPs* are used to decrease the translocation distance by changing the signs of their end genes. For more details, see [13] and [4].

##### B. Analysis of the implementation

The Algorithm 1 is a high level abstraction of the implementation of the  $1.5+\varepsilon$  algorithm that is also available at [www.mat.unb.br/~ayala/publications.html](http://www.mat.unb.br/~ayala/publications.html) as part of the whole development. The implementation uses the language C. Here, we present an overview about the running time complexity, showing that the implementation has the same complexity given in [13].

Let  $n$  be the size of the genome  $A$ . At line 1, the breakpoint graph  $G_u(A, B)$  is built using adjacency lists in time  $O(n^2)$ .



At line 2, the process of computing the 1-cycles, and at line 3 building the graph  $F_{AB}$  have time complexity  $O(n)$  each one, since it is necessary to process  $n$  genes in  $A$ .

At line 4, for computing the maximum matching graph  $M$  of  $F_{AB}$  we used the `boost` library, this well-known library is implemented in C++ and is available at <http://www.boost.org/>. Computation of the matching graph has time complexity  $O(V^2)$  with  $V$  representing the number of vertices in  $F_{AB}$ .

At line 5, finding isolated 2-cycles and 2-cycles in related components of  $M$ , has time complexity  $O(m^2)$  with  $m$  being the number of vertices of  $M$ , since it is necessary to compare if two cycles share the same grey edges in  $M$ .

At line 6, the procedure to identify the *RMSPs* has time complexity  $O(mp)$ , with  $p$  representing the number of genes of each isolated cycle.

For the lines 7 and 8, distributing appropriate signals for both 2-cycles either isolated or related, the necessary time is  $O(m^2)$ .

At line 9, removing *RMSPs* is performed in  $O(m)$ , this procedure is very simple, because reverses only the extreme genes of each *RMSP*.

At line 10, getting the signals distributed for the 2-cycles in the previous steps and assigning it to the genes of  $A$  have time complexity  $O(n)$ .

At line 11, verifying if there exist genes without signals in  $A$  have time complexity  $O(n)$ .

At line 12, distributing arbitrarily signals to genes in  $A$  have time complexity  $O(n)$ .

Thus, if one looks only to the procedure with highest complexity, that is the procedure that computes the graph  $G_u(A, B)$  and has quadratic complexity, the implementation runs as proposed in ([12], [13]).

---

**Algorithm 1:**  $1.5+\varepsilon$  approximate algorithm for UTD

---

**Data:** Genomes  $A$  (and  $B$  as identity genome)

**Result:** Genomes  $A$

- 1 Build the breakpoint graph  $G_u(A, B)$ ;
  - 2 Compute all possible 1-cycles in  $A$ ;
  - 3 Build the graph  $F_{AB}$ ;
  - 4 Compute the maximum matching graph  $M$  of  $F_{AB}$ ;
  - 5 Compute isolated 2-cycles and related components in  $M$ ;
  - 6 Build all possible *RMSPs* of isolated 2-cycles;
  - 7 Distribute appropriate signals to isolated 2-cycles;
  - 8 Distribute appropriate signals to related components;
  - 9 Remove all *RMSPs*;
  - 10 Get the signals distributed for the 2-cycles in the previous steps and put in genes of  $A$ ;
  - 11 **if** *there exists genes without signals in  $A$*  **then**
  - 12   └ Distribute arbitrarily signals to genes in  $A$ ;
- 

## V. A GENETIC ALGORITHM FOR UTD

Initially, necessary concepts about genetic algorithms are given. A GA is a searching technique used to solve optimization problems, that was introduced in 1975 by Holland in

his book "*Adaptation in natural and artificial systems*". Such technique works with the hypothesis that the genetic information of a specific population contains a possible solution. This solution, possibly is not contained in a single individual. Thus, through techniques of genetic combination, new individuals can be obtained that improve the solution to the proposed problem, and after some generations the individuals converge to a good solution [16].

In order to model a solution based in natural evolution, GAs emulate the evolutionary process that is done in the nature. The following concepts are necessary in order to understand GA:

**Individual:** an individual represents a unique solution, within a scenario of possible solutions.

**Population:** a population is a set of individuals constituting a scenario that contains a part of the search space, such population may contain potential solutions.

**Fitness:** it is used to measure how good an individual is in the population.

The process or cycle of reproduction is the central part of GAs, since during this process new individuals are created potentially with incremental quality. The reproduction cycle consist in 4 steps:

**Selection:** choose two parents to perform the reproduction. The objective is to find good individuals hoping the generation of descendants with better fitness.

**Crossover:** apply the crossover over the selected parents producing two new descendants. In the context of our problem this operation is performed by swapping the elements from a random point to the end of the string of two parents solutions.

**Mutation:** after the crossover, some of the individuals are subjected to mutation. The mutation has the roles of recovering the lost genetic material and also maintaining the genetic diversity. In the context of our problem this operation is performed by simply swapping the signs of a random element of an individual.

**Replacement:** consists in the replacement of individuals in the old population, which are the ones with lower fitness.

### A. Fitness function in the GA for UTD

The purpose of the fitness in our algorithm is to calculate the translocation distance between the signed genome  $\vec{A}$  and the identity genome, which is the genome with all its elements positive and sorted in increasing order. Thus, we can rank the best signed versions of the unsigned genome  $A$  according to this fitness. The linear algorithm proposed by Bergeron et al in [5] is used to calculate the translocation distance between two signed genomes. Originally, this linear algorithm was implemented in Java as a part of the system UniMoG [10], but we reimplemented it in the C language.

Finding an optimal solution for a given unsigned genome  $A$  is a hard task, since, the search space for such unsigned genome consists of  $2^n$  signed genomes, that are all possible signed versions of  $A$ . Such signed genomes can be sorted in linear time. Additionally, it is easy to note that solutions that solve any signed genome in the search space also solve the initial unsigned genome, and of course, that all these solutions

will require a number of translocations greater than or equal to the translocation distance of the given unsigned genome  $A$ . This fact will be used to guide the proposed GA.

### B. Description of the GA

The GA works as follows. Initially, a random population of signed genomes is generated based on the unsigned genome input. After that, for each generation the reproduction is performed as follows: Select two individuals of the population, such individuals are part of the best current solutions for which crossover and mutation operations are applied producing two new individuals. Then, the new individuals are incorporated in the current population. The GA finishes after all the generations have been completed, the number of generations depends on the size of the input genome.

The pseudo-code of our proposed standard GA is shown in Algorithm 2.

---

#### Algorithm 2: GA for Calculating UTD

---

**Input:** Unsigned Genome  $A$   
**Output:** Number of translocations to sort genome  $A$

- 1 Generate the initial population of signed genomes;
- 2 Compute fitness of the initial population;
- 3 **for**  $i = 1$  to  $Length(A)$  **do**
- 4     Perform the selection and save the best solution found;
- 5     Apply the crossover operator;
- 6     Apply the mutation operator;
- 7     Compute the fitness of the current population;
- 8     Perform replacement of the worst individuals;

---

Let  $n$  be the size of the genome  $A$ . The initial population size is defined as  $n \log n$ . Each individual in the population is generated from  $A$  in linear time, randomly assigning either a positive or negative sign to each gene. This step has time complexity of  $O(n^2 \log n)$ .

Since, for a single individual the fitness is computed in linear time, the process of computing the fitness value for all population has time complexity  $O(n^2 \log n)$ .

In the Selection step, for sorting the population by fitnesses in ascending order, we use the counting sort algorithm that runs with complexity  $O(n + n \log n)$ , with the fitness value of each individual in the interval from 1 to  $n$ , and with population size being  $n \log n$ . Thus, the complexity of this step is  $O(n \log n)$ .

In the crossover step, the best individuals classified during the selection step are chosen to be the parents. For each pair of parents, apply crossover on them by copying the elements at the right side of a random point for one individual to the other and vice versa, clearly this takes linear time ( $O(n)$ ). Thus, the running time for executing the crossover over a maximum of  $n \log n$  individuals is  $O(n^2 \log n)$ .

In the mutation step, this operator is applied to each new individual produced by the crossover. For each element of one individual, a check is made to verify whether to apply or not the mutation over a single element, this clearly takes linear time ( $O(n)$ ). The total time taken by the mutation applied over a maximum of  $n \log n$  individuals is  $O(n^2 \log n)$ .

In the replacement step, each replacement of an individual takes linear time  $O(n)$ , since we must copy all of its elements. The total time taken by the replacement of a maximum of  $n \log n$  individuals takes a total time of  $O(n^2 \log n)$ .

Finally, the genetic algorithm finishes after  $n$  generations and its total time complexity is  $O(n^3 \log n)$ .

## VI. EXPERIMENTS AND RESULTS

In order to validate the proposed GA, several tests were performed. The tests were done for randomly generated genomes. These genomes were created as follows: Generate an identity genome containing  $n$  genes and  $N$  chromosomes. Then, over this identity genome apply a fixed number of random reversals and translocations. The pseudocode of this procedure is shown in Algorithm 3.

---

#### Algorithm 3: Construction of synthetic genomes

---

**Input:** Number of genes  $n$  with  $N$  chromosomes  
**Output:** A synthetic genome  $A$

- 1 Generate an identity genome  $A$  with  $n$  genes and  $N$  chromosomes;
- 2  $j \leftarrow 0$ ;
- 3 **while**  $j \leq n$  **do**
- 4     Choose randomly a chromosome  $C$  of  $A$ ;
- 5     Select randomly an interval in  $C$ ;
- 6     Apply a reversal over this interval;
- 7     Choose randomly two chromosomes  $C$  and  $C'$  of  $A$ ;
- 8     Apply a Prefix-Prefix translocation between segments of  $C$  and  $C'$ ;
- 9      $j \leftarrow j + 1$ ;

---

In order to obtain solutions with good quality, adjustments were performed in the parameters of the genetic operators. The fine-tuning was empirically performed and provided better solutions when compared with the GA solutions without adjustments in these operators. The experiment was performed as follows: The GA was executed ten times for each genome contained in a group of hundred elements, with each group containing genomes with  $n$  genes, for  $n \in \{20, 50, 100, 150\}$ , and with 25% of each group having  $N$  chromosomes, with  $N \in \{2, 3, 4, 5\}$ . For each parameter to be adjusted its value was varied over a scenario of possible good values, and for the other parameters were fixed estimated values.

At the end of the experiment the parameters that provided the best results for the GA were taken. Those parameters are the following: single crossover point with probability of 90%, mutation probability of 2%, selection applied over 80% of the current population, and replacement applied over the 70% of the worst individuals of the current population.

Choosing these fine-tuning parameters, experiments were performed for calculating the translocation distances for the standard GA and the  $1.5+\epsilon$ -approximation algorithm. For this purpose, genomes were generated using the Algorithm 3 with  $n$  genes, for  $n \in \{10, 20, \dots, 150\}$ , and with  $N$  chromosomes, for  $N \in \{2, 3, 4, 5\}$ . For hundred genomes of length  $(n, N)$ , the average of the results for the  $1.5+\epsilon$ -approximation algorithm was calculated. The same packages of hundred genomes used in the approximation algorithm were used to calculate the

TABLE I. RESULTS OF THE  $1.5+\epsilon$ -APPROXIMATION AND THE STANDARD GA FOR 2 AND 3 CHROMOSOMES.

n	2 chromosomes		3 chromosomes	
	Average GA	Average Aprox	Average GA	Average Aprox
10	3.394	3.540	2.750	2.900
20	9.738	10.770	9.244	10.360
30	16.513	18.690	15.891	18.110
40	23.600	27.080	22.442	25.730
50	29.9810	34.580	29.662	34.170
60	37.183	42.910	36.639	42.010
70	44.907	51.840	43.230	49.930
80	51.869	59.620	50.604	58.490
90	58.213	66.960	57.715	66.620
100	66.287	76.300	65.448	75.320
110	74.534	85.940	72.940	83.710
120	80.587	92.400	80.064	91.710
130	89.164	102.020	86.838	99.590
140	96.252	110.070	94.599	108.210
150	103.510	118.380	102.106	116.350

TABLE II. RESULTS OF THE  $1.5+\epsilon$ -APPROXIMATION AND THE STANDARD GA FOR 4 AND 5 CHROMOSOMES.

n	4 chromosomes		5 chromosomes	
	Average GA	Average Aprox	Average GA	Average Aprox
10	1.670	1.740	0.980	0.980
20	8.442	9.170	7.320	7.890
30	14.861	16.710	13.650	15.330
40	21.058	23.860	20.078	22.420
50	27.545	31.450	26.334	30.020
60	34.314	39.380	32.111	36.880
70	41.488	47.640	38.935	44.430
80	47.589	54.440	45.134	51.490
90	55.020	63.110	52.609	60.310
100	61.539	70.830	58.815	67.260
110	68.687	78.460	65.027	74.450
120	75.462	86.140	72.028	82.270
130	82.452	94.020	78.792	89.680
140	89.948	102.750	86.133	98.230
150	97.686	111.080	92.506	105.330

average in the standard GA. It is worth mentioning that for each genome of length  $(n, N)$  the standard GA was executed ten times and then, the average of the ten obtained results was calculated. This average represents the result for each genome of length  $(n, N)$ .

The standard GA and the  $1.5+\epsilon$ -approximation algorithm were implemented in C language and executed in OS X platforms with Intel core I5 processors. The source code is available at [www.mat.unb.br/~ayala/publications.html](http://www.mat.unb.br/~ayala/publications.html).

The results (average translocation distances) of the experiment are shown in the Tables I and II. Also, experiments for calculating the average running time (in seconds) of 100 executions were performed for both algorithms and the results are shown in the Tables III and IV.

## VII. DISCUSSION

A few considerations are necessary before discussing the results. On the way to build synthetic genomes, instead of applying just prefix-prefix translocations between the chromosomes, we also apply reversals over the chromosomes. By including reversals it was possible to obtain harder instances of the problem. Prefix-suffix translocations are not considered because they are analogous to apply a reversal and a prefix-prefix translocations, which are already included in Algorithm 3.

TABLE III. RUNNING TIME (IN SECONDS) OF THE  $1.5+\epsilon$ -APPROXIMATION AND THE STANDARD GA FOR 2 AND 3 CHROMOSOMES

n	2 chromosomes		3 chromosomes	
	Average GA	Average Aprox	Average GA	Average Aprox
10	0.008	0.010	0.009	0.010
20	0.030	0.010	0.032	0.011
30	0.082	0.010	0.091	0.012
40	0.184	0.009	0.199	0.011
50	0.356	0.010	0.377	0.010
60	0.605	0.010	0.646	0.019
70	0.962	0.010	1.029	0.010
80	1.427	0.010	1.535	0.011
90	2.076	0.010	2.186	0.010
100	2.877	0.010	3.016	0.010
110	3.839	0.010	4.032	0.011
120	5.011	0.011	5.246	0.010
130	6.498	0.011	6.811	0.011
140	8.123	0.011	8.519	0.011
150	10.071	0.011	10.489	0.011

TABLE IV. RUNNING TIME (IN SECONDS) OF THE  $1.5+\epsilon$ -APPROXIMATION AND THE STANDARD GA FOR 4 AND 5 CHROMOSOMES

n	4 chromosomes		5 chromosomes	
	Average GA	Average Aprox	Average GA	Average Aprox
10	0.011	0.013	0.012	0.012
20	0.037	0.015	0.042	0.011
30	0.101	0.014	0.110	0.011
40	0.217	0.011	0.234	0.011
50	0.405	0.010	0.434	0.010
60	0.685	0.009	0.737	0.010
70	1.082	0.010	1.134	0.010
80	1.626	0.010	1.687	0.010
90	2.322	0.010	2.415	0.010
100	3.170	0.010	3.322	0.010
110	4.258	0.010	4.413	0.010
120	5.479	0.011	5.756	0.010
130	7.091	0.011	7.384	0.011
140	8.899	0.011	9.273	0.011
150	10.975	0.011	11.421	0.011

It is important to emphasize that the algorithm [5] used as fitness function had already been implemented by their authors as a contribution of Jens Stoye. The source code implemented in Java was made available and then translated to C. Also, we performed several tests to validate the translocation distance. For the validation, the algorithm proposed in [4] (as corrected by Bergeron in [5]) was used, which provides the optimal sequence of translocations necessary to transform a genome into another.

There is a little variation in the running time of the  $1.5+\epsilon$  approximation algorithm even for genomes of length 150. This is because the steps of the algorithm are relatively simple and also since, the solutions of the approximation algorithm are based in the calculation of 2-cycles and the inputs generated randomly have a few number of 2-cycles. So, the execution of the algorithm is always fast.

For the standard GA, we can observe that when the number of chromosomes and genes are incremented, the running time grows rapidly. This is because the size of the population depends on the number of genes. Although this, it is necessary to stress here that the size of the population is not proportional for the search space for genomes of different size as usual in combinatorics of permutations ( $n \log n$  versus  $n!$ ).

From the experiments, one can conclude that the standard GA compute better results on average than those obtained

by the  $1.5+\varepsilon$ -approximation algorithm. It can be observed in Tables I and II that for permutations of length greater than or equal to 50, the standard GA has better solutions in approximately 12%.

As can be seen in the Tables III and IV the running time of the standard GA is, as expected, greater when compared with the running time of the approximation algorithm and this difference is higher for larger inputs.

## VIII. CONCLUSIONS AND FUTURE WORK

In the search for good solutions for the  $\mathcal{NP}$ -hard problem of translocation distance for unsigned genomes, a standard genetic algorithm was proposed in this paper. This standard GA acts on a population of signed genomes generated from an unsigned genome (the input), and after each generation the population evolves to the signed genomes with the best translocation distance. Indeed, the distinguished feature of our GA is using as fitness function the (linearly computable) translocation distance of signed permutations.

The experiments showed that results obtained by the standard GA outperform the results obtained by the  $1.5+\varepsilon$ -approximation algorithm. With respect to the running time, the  $1.5+\varepsilon$ -approximation algorithm as expected is faster when compared with the standard GA, however this difference is tolerable, since, the experiments with the standard GA have running time of approximately 10 seconds for genomes with 150 genes.

As an immediate further step we are planning experiments with data generated from real genomes, which can be obtained from the biological sequence database *GeneBank*. This data would be generated by assigning an integer number to each gene of a real genome; these integer numbers are mapped from an identity genome with the same genes. Also, as future work we are planning to improve our standard GA by including other interesting heuristics as done for other GA approaches to deal with reversal distance. This will include memetic GA approaches as in [17] and parallel GA approaches as in [18] with the aim of improving the quality of the solutions and reducing the running time. Also it would be of great interest performing experiments with opposition based learning with the aim of exploring the search space using opposite solutions [19].

## ACKNOWLEDGMENT

The authors would like to thank the Brazilian agencies CAPES and CNPq for funding this work with scholarships and a universal grant.

## REFERENCES

- [1] V. Bafna and P. A. Pevzner, "Sorting by transpositions," *SIAM Journal on Discrete Mathematics*, vol. 11, no. 2, pp. 224–240, 1998.
- [2] S. Hannenhalli and P. A. Pevzner, "Transforming cabbage into turnip: polynomial algorithm for sorting signed permutations by reversals," *Journal of the ACM (JACM)*, vol. 46, no. 1, pp. 1–27, 1999.
- [3] D. A. Bader, B. M. Moret, and M. Yan, "A linear-time algorithm for computing inversion distance between signed permutations with an experimental study," *Journal of Computational Biology*, vol. 8, no. 5, pp. 483–491, 2001.

- [4] S. Hannenhalli, "Polynomial-time algorithm for computing translocation distance between genomes," *Discrete Applied Mathematics*, vol. 71, no. 1, pp. 137–151, 1996.
- [5] A. Bergeron, J. Mixtacki, and J. Stoye, "On sorting by translocations," *Journal of Computational Biology*, vol. 13, no. 2, pp. 567–578, 2006.
- [6] D. Zhu and L. Wang, "On the complexity of unsigned translocation distance," *Theoretical computer science*, vol. 352, no. 1, pp. 322–328, 2006.
- [7] A. Caprara, "Sorting by reversals is difficult," in *Proceedings of the first annual international conference on Computational molecular biology*, ser. RECOMB '97. New York, NY, USA: ACM, 1997, pp. 75–83. [Online]. Available: <http://doi.acm.org/10.1145/267521.267531>
- [8] I. Holyer, "The np-completeness of some edge-partition problems," *SIAM Journal on Computing*, vol. 10, no. 4, pp. 713–717, 1981.
- [9] L. Wang, D. Zhu, X. Liu, and S. Ma, "An  $o(n^2)$  algorithm for signed translocation," *Journal of Computer and System Sciences*, vol. 70, no. 3, pp. 284–299, 2005.
- [10] R. Hilker, C. Sickinger, C. N. Pedersen, and J. Stoye, "Unimog—a unifying framework for genomic distance calculation and sorting based on dcj," *Bioinformatics*, vol. 28, no. 19, pp. 2509–2511, 2012.
- [11] J. D. Kececioğlu and R. Ravi, "Of mice and men: algorithms for evolutionary distances between genomes with translocation," in *Proceedings of the sixth annual ACM-SIAM symposium on Discrete algorithms*. Society for Industrial and Applied Mathematics, 1995, pp. 604–613.
- [12] Y. Cui, L. Wang, and D. Zhu, "A  $1.75$ -approximation algorithm for unsigned translocation distance," *Journal of Computer and System Sciences*, vol. 73, no. 7, pp. 1045–1059, 2007.
- [13] Y. Cui, L. Wang, D. Zhu, and X. Liu, "A  $(1.5+\varepsilon)$ -approximation algorithm for unsigned translocation distance," *IEEE/ACM Transactions on Computational Biology and Bioinformatics (TCBB)*, vol. 5, no. 1, pp. 56–66, 2008.
- [14] H. Jiang, L. Wang, B. Zhu, and D. Zhu, "A  $(1.408+\varepsilon)$ -approximation algorithm for sorting unsigned genomes by reciprocal translocations," in *Frontiers in Algorithmics*. Springer, 2014, pp. 128–140.
- [15] R. M. Karp, *Reducibility among combinatorial problems*. Springer, 1972.
- [16] S. Sivanandam and S. Deepa, *Introduction to genetic algorithms*. Springer, 2007. [Online]. Available: <http://books.google.com.br/books?id=wonrLjj2GagC>
- [17] J. L. Soncco-Álvarez and M. Ayala-Rincón, "Memetic algorithm for sorting unsigned permutations by reversals," in *Evolutionary Computation (CEC), 2014 IEEE Congress on*. IEEE, 2014, pp. 2770–2777.
- [18] J. L. Soncco-Álvarez, G. M. Almeida, J. Becker, and M. Ayala-Rincón, "Parallelization of genetic algorithms for sorting permutations by reversals over biological data," *International Journal of Hybrid Intelligent Systems*, vol. 12, no. 1, pp. 53–64, 2015.
- [19] Q. Xu, L. Wang, N. Wang, X. Hei, and L. Zhao, "A review of opposition-based learning from 2005 to 2012," *Engineering Applications of Artificial Intelligence*, vol. 29, pp. 1–12, Mar. 2014. [Online]. Available: <http://linkinghub.elsevier.com/retrieve/pii/S0952197613002388>