

A sparse dynamic programming algorithm for alignment with non-overlapping inversions

Alair Pereira do Lago Ilya Muchnik
Casimir Kulikowski

September 26, 2003

Abstract

Alignment of sequences is widely used for biological sequence comparisons, and only biological events like mutations, insertions and deletions are considered. Other biological events like inversions are not automatically detected by the usual alignment algorithms, thus some alternative approaches have been tried in order to include inversions or other kind of rearrangements.

Despite many important results in the last decade, the complexity of the problem of alignment with inversions is still unknown. In 1992, Schöniger and Waterman proposed the simplification hypothesis that the inversions do not overlap. They also presented an $O(n^6)$ exact solution for the alignment with non-overlapping inversions problem, and introduced a heuristic for it that brings the running-time complexity down.

The present paper gives two exact algorithms for the simplified problem. We give a quite simple dynamic program with $O(n^4)$ -time and $O(n^2)$ -space complexity for alignments with non-overlapping inversions and exhibit a sparse and exact implementation version of this procedure that uses much less resources for some applications with real data.

1 Introduction

In evolution history, some biological events introduce changes in the DNA sequences. Some typical biological events include *mutations*, in which a nu-

cleotide is substituted by another one, *deletions* and *insertions* of nucleotides. Hence, any sequence comparison must take into consideration the possibility of these events if it is expected to identify high similarity between two sequences. Typical alignment procedures try to identify which parts do not change and where these biological events are, after exhibiting one best alignment according to some optimization criteria.

For instance, in Figure 1, we see two alignments of *actagatcagtca* against *attgaatcgacta*. From left to right, one can detect a deletion of *c*, a mutation

actaga-tcagtc-a	actaga-tcagtca

a-ttgaatc-gacta	a-ttgaatcgacta

Figure 1: Two alignments for *actagatcagtca* and *attgaatcgacta*

from *a* to *t* and an insertion of *a* in both alignments. At the rightmost part, the first alignment report a deletion of *a*, a mutation from *t* to *a* and an insertion of *t*. In contrast, the second alignment highlights the inversion of *agtc* to its reverse complement *gact*. Since common aligners do not take *inversions* into consideration, they would report the first alignment. Figure 2 shows how an inversion can occur and why one segment is substituted by its reverse complement sequence. Alignments can be associated with a set of edit operations that transform one sequence to the other. Usually, the only edit operations that are considered are the *substitution* (mutation) of one symbol by another one, the *insertion* of one symbol and *deletion* of one symbol. If costs are associated with each operation, there is a classic $O(n^2)$ dynamic program that computes a set of edit operations with minimal total cost and exhibit the associated alignment, which has good quality and high likelihood for realistic costs.

Consider three new possible edit operations:

- the *2-reversion*, which reverses the order of *two consecutive symbols*;
- the *reversion* operation, which reverses the order of *any segment* of symbols instead of a segment of length 2;
- the *inversion* operation, which substitutes any segment by its *reverse complement* sequence. The inversion operation is the operation that is interesting in molecular biology. (See Figure 2.)

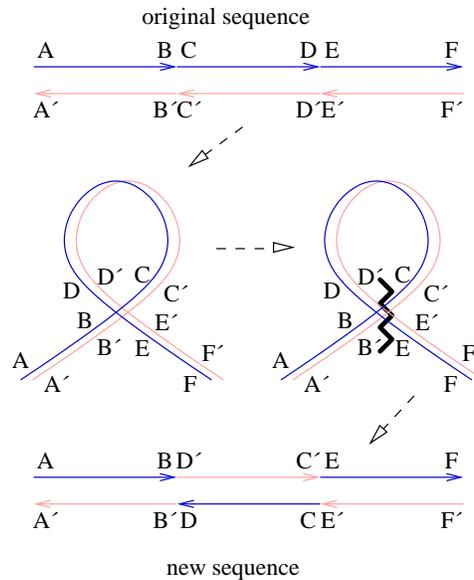


Figure 2: Example of DNA inversion

Associated with any of these three operations, we can define new alignment problems. For instance, given two sequences and fixed costs for each kind of edit operation, the *alignment with inversions* problem is an optimization problem that queries the minimal total cost of an edit operations set that transforms one sequence to the other. Moreover, one may also be interested in the exhibition of its correspondent alignment and/or edit operations. Similarly, we can define the *alignment with 2-reversions* and the *alignment with reversions* problems.

In 1975, Wagner [1] studied the alignment with 2-reversions problem and proved that it admits a polynomial solution if the cost of 2-reversion is null. On the other hand, he also proved that the obtainment of an optimal solution is *NP-hard*, if any operation has a constant positive cost.

To the best of our knowledge, the computational complexities of alignment with reversions and alignment with inversions problems are unknown.

In order to deal with alignments with inversions, three main approaches have been considered through the years:

- non-overlapping inversions;
- sorting unsigned permutations by reversals and;

- sorting signed permutations by reversals.

Before we proceed with the results of this paper, we will give a brief summary of these three approaches.

In 1992, Schöniger and Waterman [2] introduced a *simplification hypothesis*: *all regions involving the inversions do not overlap*. This led to the *alignment with non-overlapping inversions* problem. They presented a $O(n^6)$ solution for this problem and also introduced a *heuristic* for it that reduced the running-time. This heuristic uses the algorithm by Waterman and Eggert [3] that reports the K best non mutual intersecting local alignments in order to reduce the running time to something between $O(n^2)$ and $O(n^4)$, depending on the data.

Another approach has been tried in order to study inversions. This approach applies well for alignment of *sequences of genes* and has been very used with mitochondrial genomes. It does not apply for sequences of nucleotides nor for sequences of aminoacids because *no repetitions* of symbols are allowed. (Repeated genes and paralogs are not allowed.) Moreover, *no insertion* and *no deletion* are considered and the only admitted operation is the reversal, where a *reversal* is defined to transform a sequence like $(1, 2, 3, 4, 5)$ into $(1, 4, 3, 2, 5)$.

The problem, also called *sorting unsigned permutations by reversals*, means the computation of the edit distance of two permutations when only the reversal operation is allowed. In this case, the data are two permutations of $(1, 2, 3, \dots, n)$, where n is the number of genes. Kececioglu and Sankoff [4] gave a quadratic 2-approximation algorithm in 1995 and Christie [5] gave a $3/2$ -approximation algorithm in 1998. Caprara [6] proved in 1999 that this problem is in fact NP-hard.

Another interesting approach is the problem called *sorting signed permutations by reversals*. This is the same problem as sorting unsigned permutations by reversals up to the fact that signals are also attributed to a gene and a reversal also flips its signal. For instance, one reversal could transform $(1, 2, 3, 4, 5)$ to $(1, -4, -3, -2, 5)$. This signal is usually associated with the direction of the gene (which DNA strand it belongs to).

Hannenhalli and Pevzner [7, 8] gave the first polynomial algorithm for the problem in 1995 and started a sequence of papers based on this approach. Hannenhalli and Pevzner's algorithm was $O(n^4)$ and it was improved to $O(n^2)$ by Kaplan, Shamir and Tarjan [9, 10] in 1997. In 2001, Bader, Moret and Yan [11] gave an algorithm that computes the minimal number of reversals

distance in $O(n)$ (finding the sequence of reversals still requires $O(n^2)$). These studies have been applied for phylogenetic reconstruction studies.

In 2000, El-Mabrouk [12, 13] studied the inclusion of two operations: insertions and deletions of gene segments. She obtained partial results and gave an exact polynomial solution for one case and a polynomial heuristic with a polynomial tester for optimality in the other case. Repeated symbols are still not allowed. In 2002, El-Mabrouk [14] also obtained some partial results on considering reversals and duplications.

This paper gives two exact algorithms for the alignment with non-overlapping inversions problem, which is the first approach. Algorithm 1 is a $O(n^4)$ solution that uses $O(n)$ space and Algorithm 2 is a sparse dynamic implementation of it that reduces the resources usage if there are $o(n^2)$ matches. This is often expected if the cardinality of the alphabet is large and it is true for the kind of application we have in mind, where the letters are DNA fragments of fixed length.

2 Basic Definitions

Let A be any *alphabet*, a set of *letters*. Any finite sequence on A is also called a *word on A* or simply a *word* if the alphabet is clear. Let A^* be the set of all words on A , including the *empty word* denoted by 1 . We identify words of length 1 to the letters they contain. The concatenation \cdot of words is an associative operation defined over A^* and it will be often omitted. Let $w = w_1w_2 \dots w_k$ be a word. We denote by $|w|$ the *length k* of w . We will also denote w_i , the i -th letter of w , by $w[i]$. Let $x, y, z \in A^*$. We denote by xA^* the set $\{xy \mid y \in A^*\}$, we denote by A^*x the set $\{yx \mid y \in A^*\}$ and we denote by A^*xA^* the set $\{yxz \mid y, z \in A^*\}$. We say that x is a *prefix of w* if $w \in xA^*$, we say that x is a *suffix of w* if $w \in A^*x$ and we say that x is a *factor of w* if $w \in A^*xA^*$. For $1 \leq i \leq j \leq k$, the factor $w_iw_{i+1} \dots w_j$ of w is also represented by $w[i..j]$.

Let $\bar{}$, also called *inversion*, be any operation on A^* that satisfies the following properties:

1. $\bar{\bar{a}} \in A, \forall a \in A$
2. $\overline{x \cdot y} = \bar{y} \cdot \bar{x}, \forall x, y \in A^*$

Notice that the inversion operation on A^* is defined by its values on the letters of A . For instance, let $A = \{a, c, t, g\}$ and let $s \in A^*$ be any DNA sequence. If the inversion is defined by $\bar{a} = a, \bar{t} = t, \bar{c} = c$ and $\bar{g} = g$, it maps s to its *reverse sequence*. On the other hand, if the inversion is defined by $\bar{a} = t, \bar{t} = a, \bar{c} = g$ and $\bar{g} = c$, it maps s to its *reverse complement sequence*. Last case is the interesting case for DNA sequences in molecular biology.

Let $\omega : A \times A \rightarrow \mathbb{R} \cup \{-\infty\}$ be any *weight function*. We say that $(a, b) \in A \times A$ is a *match* if $\omega(a, b) \neq -\infty$.

Let $s = s_1s_2 \cdots s_k$ and $t = t_1t_2 \cdots t_{k'}$ be two words. We define the *matching graph of s and t* as the weighted and colored bipartite graph $G = G(s, t, \omega, -) = (V, E)$. The vertex set is a set of symbols $V = \{s_1, \dots, s_k, t_1, \dots, t_{k'}\}$ and has size $|s| + |t|$. (We do not identify two symbols s_i and t_j even in the case the letters s_i and t_j are the same letter in A .) The edge set E is a double copy of $K_{|s|, |t|}$: for any pair of vertices s_i and t_j we link them with one *blue/dark-gray edge of weight $\omega(s_i, t_j)$* and one *red/light-gray edge of weight $\omega(s_i, \bar{t}_j)$* . (In fact, edges with weight $-\infty$ will not be used for our purposes and could be deleted.) An edge with weight that is not $-\infty$ is called a *direct match* if it is blue and it is called *inverted match* if it is red. In Figure 3 we have an example of a matching graph. In this figure, as in others, we do not draw edges with weight $-\infty$.

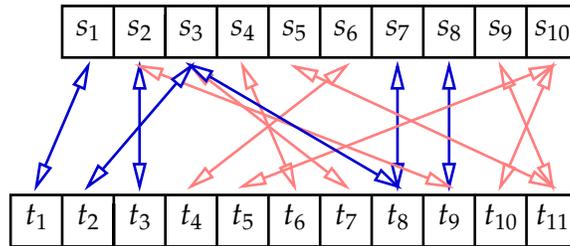


Figure 3: Example of a matching graph

Given $u \in V^*$ a nonempty factor of $s_1s_2 \cdots s_k$ and $v \in V^*$ a nonempty factor of $t_1t_2 \cdots t_{k'}$, we call $B = (u, v)$ a *block*. Given a block $B = (u, v)$, there exist integers $1 \leq i' \leq i \leq |s|$ and $1 \leq j' \leq j \leq |t|$ such that $u = s[i' \dots i]$ and $t = t[j' \dots j]$. Hence, we can associate $([i' \dots i] \times [j' \dots j])$, the *rectangle associated with the block B* . In Figure 4 we have an incidence matrix of the matching graph of Figure 3 and only matches are shown. Cells at position (s_i, t_j) are colored according to the color of the edge linking these vertices. (We could use purple if we had both matches for the same cell.) We can also see

four rectangles corresponding to the blocks $(s_1s_2, t_1t_2t_3)$, $(s_3s_4s_5s_6, t_4t_5t_6t_7)$, (s_7s_8, t_8t_9) and $(s_9s_{10}, t_{10}t_{11})$. Cells inside a rectangle correspond to edges with vertices in the factors that form the respective block.

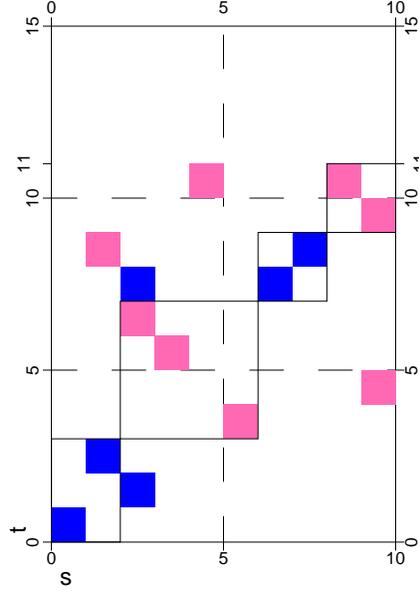


Figure 4: Bicoloured incidence matrix

Consider an edge that links s_i to t_j and an edge that links $s_{i'}$ to $t_{j'}$. We say that they *cross* each other if $(i - i')(j - j') < 0$, we say that they *touch* each other if $(i - i')(j - j') = 0$ and we say that they are *parallel* if $(i - i')(j - j') > 0$. Let $M \subseteq E$ be any set of edges in a matching graph. Recall that M is called a *matching* if any two edges of M do not touch each other. Moreover, M is called a *direct matching* if it has only direct matches and any two of them are parallel. Furthermore, M is called an *inverted matching* if it has only inverted matches and any two of them cross each other. The *restriction of M to a block $B = (s[i'..i], t[j'..j])$* is the submatching of all edges of M with vertices in $s[i'..i]$ and $t[j'..j]$. Finally, M is called a *blockwise inverted matching*, or simply a *bimatching*, if, considering words in V^* ,

- there is $l \geq 1$;
- there are l blocks $B_i = (u_i, v_i)$ such that $s = s_1s_2 \dots s_k = u_1u_2 \dots u_l$ and $t = t_1t_2 \dots t_{k'} = v_1v_2 \dots v_l$;

- for any $i \in \{1, \dots, l\}$, the restriction M_i of M to the block B_i is either a direct matching or an inverted matching;
- $M = \cup_{i=1}^l M_i$.

Given a bimatching M , the number of blocks l is not unique since any direct submatching M_i can also be split as a union of smaller direct submatchings. However, the number of blocks such that the corresponding restriction M_i is an inverted matching does not change. This is the *number of inversions of M* , $\iota(M)$. Given a bimatching M , the smallest possible value for l is called the *number of blocks of M* .

One can notice that in Figure 5 we have a bimatching M with four blocks which are those associated with the blocks of Figure 4. There are $\iota(M) = 2$ maximal inverted submatchings.

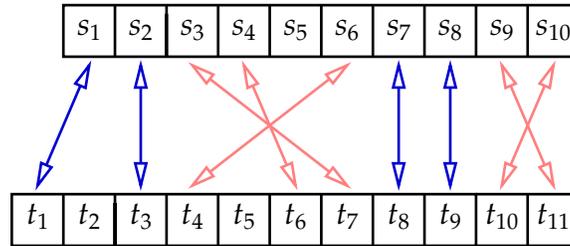


Figure 5: Example of a bimatching with four blocks

3 Two Algorithms for Optimal Bimatching

Given a bimatching M , one can easily deduce an alignment with non-overlapping inversions associated with it. One could naturally define the weight of a matching in a matching graph to be the sum of the weights of its edges. However, it is quite common for alignments the establishment of penalties for biological events like mutations, insertions or deletions. In order to be more general, we attribute a *inversion penalty* $I \geq 0$ for every inverted submatching M_i . Hence, we define the *weight of a bimatching M* as

$$\omega(M) = \sum_{e \in M} \omega(e) - \iota(M)I.$$

Hence, the following problem is the optimization problem we are interested in solving.

Problem 1 Given a matching graph $G = G(s, t, \omega, -)$, we want to compute the maximal weight $\omega(M)$ for all possible bimatings M . As usual, we may also be interested in a bimatting of maximal weight.

Such a bimatting M^* of maximal weight is called an *optimal bimatting* of s and t and its weight $\omega(M^*)$ is denoted by $\text{BIM}(s, t)$. The weight of an *optimal direct matching* is denoted by $\text{DM}(s, t)$.

Next lemma leads us to the dynamic programming algorithm given by Algorithm 1.

Lemma 2 Let A be an alphabet, let $a, b \in A$ be any letters, let $u, v \in A^*$ be any words on A , let us give an inversion operation on $-$ on A^* and a weight function ω and let I be an inversion penalty. Hence, the following affirmatives are true:

1. $\text{BIM}(1, v) = \text{BIM}(u, 1) = 0;$

2.
$$\text{BIM}(ua, vb) = \max \left\{ \begin{array}{l} \text{BIM}(u, v) + \omega(a, b) \\ \text{BIM}(ua, v) \\ \text{BIM}(u, vb) \\ B - I \end{array} \right\}$$

where $B = \max\{\text{BIM}(u_1, v_1) + \text{DM}(u_2, \bar{v}_2) \mid ua = u_1u_2, vb = v_1v_2, u_2v_2 \neq 1\}$.

(All the proofs are omitted in this paper due to space constraints.)

Algorithm 1 computes the maximal weight of a bimatting $\text{BIM}(s, t)$ with time complexity $|s|^2|t|^2$ and space complexity $|s||t|$. As usual, by tracking any maximality choice, one can obtain also an optimal bimatting. The numbers present in positions (s_i, t_j) of the incidence matrix of Figure 6 show the corresponding value $B[i, j]$ computed in Algorithm 1 if all matches have weight 1. Hence, one can obtain the bimatting of Figure 5 as the optimal bimatting for the matching graph of Figure 3.

In Algorithm 2, we improve the computational resources usage required by Algorithm 1 when there are $o(n^2)$ matches. We use techniques that appeared in works by Hunt and Szymanski [15] in 1977 for the computation of the length of the longest common subsequence (LCS) (A good survey can be found in [16]). The name *Sparse Dynamic Programming* was adopted by Eppstein et al. [17, 18] in a well known work published in 1992. The main idea behind these techniques is that only cells associated with a match are “visited”, and this is done here.

Algorithm 1 A $O(n^4)$ -time and $O(n^2)$ -space algorithm for BIM

BIM(s, t)

```

1  ▷ Compute the table  $B[i, j] = \text{BIM}(s[1..i], t[1..j])$ 
2  Let  $B[i, j]$  be 0 for  $i = 0$  or  $j = 0$ 
3  for  $i$  from 1 to  $|s|$  do
4      for  $j$  from  $|t|$  downto 1 do
5           $B[i, j] \leftarrow -\infty$ 
6          for  $j'$  from 1 to  $|t|$  do
7               $B[i, j] \leftarrow \max(B[i, j], B[i-1, j], B[i, j-1])$ 
8               $B[i, j] \leftarrow \max(B[i, j], B[i-1, j-1] + \omega(s[i], t[j]))$ 
9              ▷ Compute  $L[i', j'] = \text{DM}(s[i'..i], t[j'..j])$  and set  $B[i, j']$ 
10             Let  $L[i', j']$  be 0 for  $i' = i + 1$  or  $j' = j - 1$ 
11             for  $j'$  from  $j$  to  $|t|$  do
12                 for  $i'$  from  $i$  downto 1 do
13                      $L[i', j'] \leftarrow \max(L[i', j'-1], L[i'+1, j'])$ 
14                      $L[i', j'] \leftarrow \max(L[i'+1, j'-1] + \omega(s[i'], t[j']), L[i', j'])$ 
15                      $B[i, j'] \leftarrow \max(B[i, j'], L[i', j'] + B[i'-1, j-1] - I)$ 
16  return  $B$ 

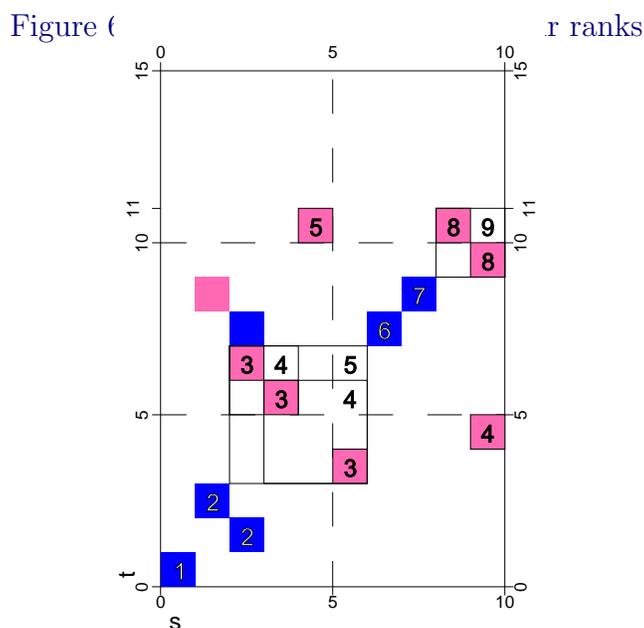
```

Before we proceed, we need some few more definitions. We say that a rectangle $([i..j] \times [k..l])$ *dominates* the rectangle $([i'..j'] \times [k'..l'])$ if $j \leq j'$ and $l \leq l'$. The pair (j, l) is called the *right upper corner* of $([i..j] \times [k..l])$ and the domination relation depends only on the right upper corners. We will give directions signals for rectangles according to whether we admit a direct or an inverted matching as the restriction of a bmatching to the corresponding block. We define the *rank* of a signed rectangle $([i..j] \times [k..l])$ to be the maximal weight of bmatchings that are restricted to the block $(s[a..j], t[c..l])$ and admit a block decomposition where the last block is $(s[i..j], t[k..l])$ and the restriction to it has the right direction. We say that a signed rectangle is *dominant* if any other signed rectangle that dominates it has either a smaller rank or the same right upper corner. The dominance relation is a partial quasi-order where the equivalence classes have always the same rank and the same right upper corner.

The only considered signed rectangles in Algorithm 2 are all direct rectangles $([i..i] \times [j..j])$ such that (i, j) is a direct match and all inverted rectangles $([i..j] \times [k..l])$ such that both (j, k) and (i, l) are inverted matches.

We store only dominant signed rectangles (one for every equivalence class). This naturally gives us the rank of their right upper corner positions. (See function `UPDATE` in Algorithm 3.) The ranks of other positions are “automatically propagated” from dominant positions. (See function `RANK` in Algorithm 3.)

Figure 6 shows dominant rectangles with their corresponding ranks in the right upper corners as computed in Algorithm 2 (all matches have weight 1). Figure 7 shows an example of possible data that have been applied to the algorithm `BIM`. The DNA sequences correspond to two syntenic regions from two bacteria. These sequences were splitted in fragments of length 100 in order to form the alphabet. A match between two fragments is assumed if the alignment score is above an adequately chosen threshold.



4 Conclusion and open problems

We gave new algorithms for the alignment with non-overlapping inversions problem, improving the time complexity of an exact solution from $O(n^6)$ to $O(n^4)$ in Algorithm 1. In Algorithm 2, we also gave a sparse dynamic programming implementation that gives the exact solution and can speed up

Algorithm 2 A $O(m \log k + (m')^2(\log k' + \log^2 k) + n)$ -time and $O(m + m' + n + k + k')$ -space algorithm for BIM, where $m = O(n^2)$ is the number of direct matches, $m' = O(n^2)$ is the number of inverted matches, $k = O(n^2)$ is the number of dominant rectangles and $k' = O(m')$ is the number of dominant direct (former inverted) matches in the computation of $\text{DM}(t, \bar{s})$

$\text{BIM}(s[a..b], t[c..d], M, \omega)$

```

1  ▷ Compute  $R$ , an AVL-tree of AVL-trees  $R_r$  of dominant rectangles
2  ▷ Any AVL-tree  $R[r] = R_r$  has only rectangles of rank  $r$ , for  $r \in \mathbb{R}$ 
3  ▷ In  $R[r]$ , rectangles  $([i'..i] \times [j..j'])$  are ordered by  $i$ 
4  ▷ In  $R$ , AVL-trees  $R[r]$  are ordered by the rank  $r$ 
5  ▷  $\text{RANK}(R, i, j) = \text{BIM}(s[a..i], t[c..j])$ 
6  ▷  $M$ : list of matches with signals  $(i, j, d) \in [a..b] \times [c..d] \times \{+1, -1\}$ 
7  ▷  $M$  is ordered by  $i$ , then by  $j$ , then by  $d$ 
8   $\delta \leftarrow$  the reflection-rotation defined by  $\delta(i, j, d) = (j, |s| + 1 - i, -d)$ 
9   $M' \leftarrow \delta(M \cap \mathbb{Z} \times \mathbb{Z} \times \{-1\})$  ▷ All inverted matches, mapped by a  $\delta$ 
10 Sort all  $(i, j, d) \in M'$  by  $i$ , then by  $j$ 
11  $R \leftarrow \emptyset$ 
12  $R[0] \leftarrow R_0 \leftarrow ([0..0] \times [0..0])$ 
13  $R[-\infty] \leftarrow \emptyset$ 
14  $R[+\infty] \leftarrow \emptyset$ 
15 for  $(i, j, d) \in M$  do
16     if  $d = +1$  then
17          $r \leftarrow \omega(s[i], t[j]) + \text{RANK}(R, i - 1, j - 1)$ 
18         if  $r \geq \text{RANK}(R, i, j)$  then
19              $\text{UPDATE}(R, r, ([i..i] \times [j..j]), +1)$ 
20     else  $M'' \leftarrow M' \cap [j..d] \times [|s| + 1 - i..|s| + 1 - a] \times \{+1\}$ 
21         ▷  $[j..d] \times [|s| + 1 - i..|s| + 1 - a] \times \{+1\} = \delta([a..i] \times [j..d] \times \{-1\})$ 
22          $\varphi \leftarrow$  function defined by  $\varphi(f, g) = \omega(\bar{g}, \bar{f})$ 
23          $B \leftarrow \text{BIM}(t[j..d], \bar{s}[|s| + 1 - i..|s| + 1 - a], M'', \varphi)$ 
24         for  $\delta(i', j', -1) \in M''$  do
25             ▷  $\text{DM}(s[i'..i], \overline{t[j..j']}) + \text{BIM}(s[a..i' - 1], t[c..j - 1]) - I$ 
26              $r \leftarrow \text{RANK}(B, j', |s| + 1 - i') + \text{RANK}(R, i' - 1, j - 1) - I$ 
27             if  $r \geq \text{RANK}(R, i, j')$  then
28                  $\text{UPDATE}(R, r, ([i'..i] \times [j..j']), -1)$ 
29 return  $R$ 

```

Algorithm 3 Functions RANK and UPDATE used in Algorithm 2

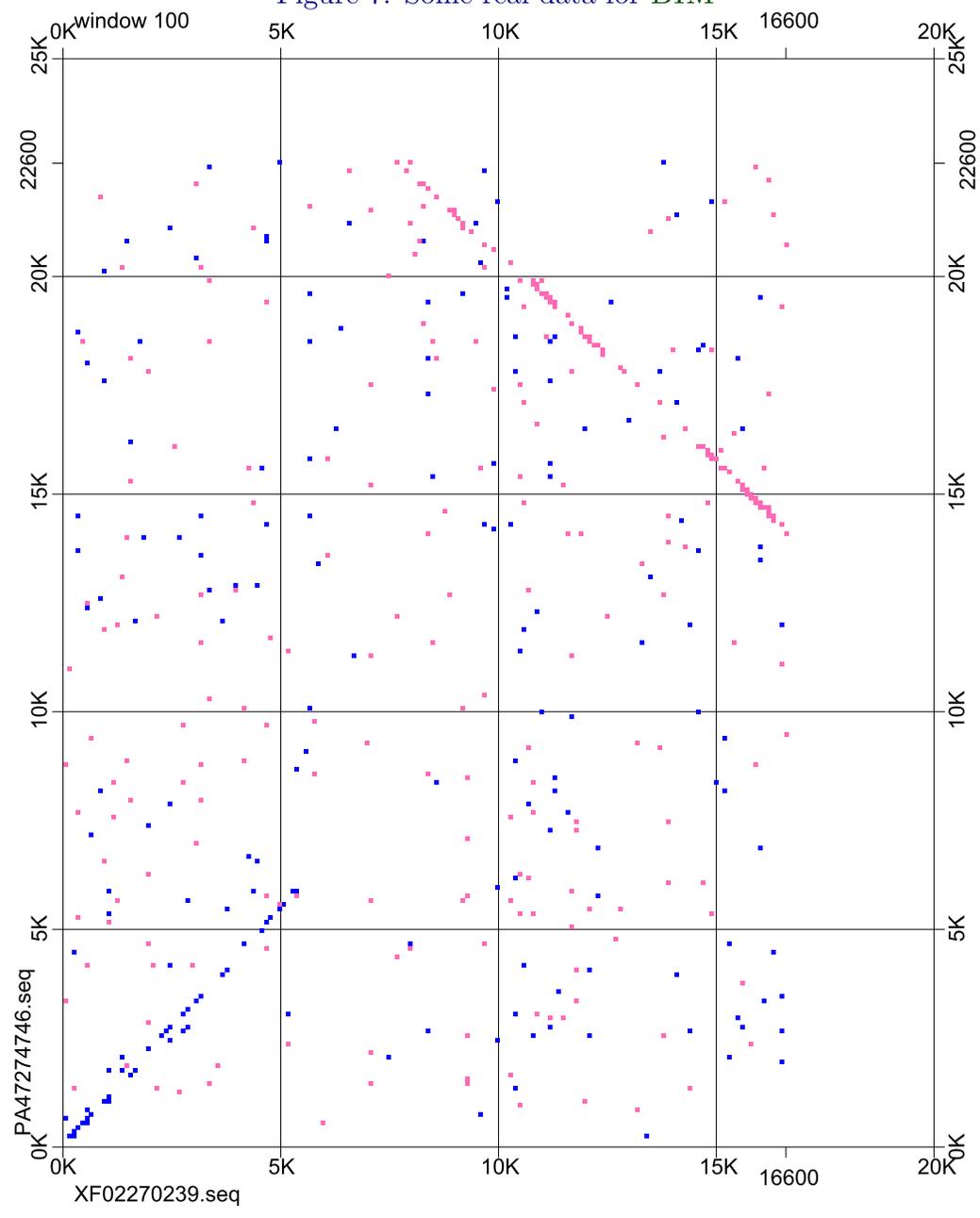
RANK(R, i, j)

```
1  ▷ Return the rank of position  $(i, j)$ 
2   $a \leftarrow -\infty$ 
3   $b \leftarrow +\infty$ 
4   $m \leftarrow$  the rank of the root of  $R$ 
5  while  $m \neq a$  and  $m \neq b$  do
6      if  $\exists([x..y] \times [u..v]) \in R[m]$  that dominates  $([i..i] \times [j..j])$  then
7           $a \leftarrow m$ 
8           $m \leftarrow$  right child of  $m$ 
9      else  $b \leftarrow m$ 
10      $m \leftarrow$  left child of  $m$ 
11 return  $a$ 
```

UPDATE($R, r, ([i..i'] \times [j..j']), d$)

```
1  ▷ Insert the rectangle  $([i..i'] \times [j..j'])$  of rank  $r$  in  $R$ 
2   $r' \leftarrow r$ 
3  do for each  $([x..y] \times [u..v]) \in R[r']$  dominated by  $([i..i'] \times [j..j'])$  do
4      Remove  $([x..y] \times [u..v])$  from  $R[r']$ 
5       $r' \leftarrow$  rank that is predecessor of  $r'$  in  $R$ 
6  while some rectangle  $([x..y] \times [u..v])$  was removed from  $R[r']$ 
7  Insert  $([i..i'] \times [j..j'])$  in  $R[r]$   ▷ We had  $r \geq \text{RANK}(R, i', j')$ 
```

Figure 7: Some real data for BIM



even further if we have $o(n^2)$ matches. This is quite common for applications with large alphabets. We can also modify the algorithms in order to deal with gap penalties or local alignments with non-overlapping inversions.

Let $\text{LCS}(u, v)$ denote the length of the longest common subsequence of u and v . Motivated by these algorithms, one of the authors recently [20] proposed the following open problem: given two words s and t of length n , one can pre-process both words in such a way that any query $\text{LCS}(u, v)$ can be answered in time $O(1)$, for u a factor of s and v a factor of t . This pre-processing can be done in $O(n^4)$. Can we do it in $O(n^2)$? In $O(n^3)$? Although a little stronger, one can think of a version with $\text{DM}(u, v)$ queries instead of $\text{LCS}(u, v)$ queries.

As far as we know, alignment with non-restricted inversions is open.

References

- [1] Wagner, R.: On the complexity of the extended string-to-string correction problem. In: Seventh ACM Symposium on the Theory of Computation, Association for Computing Machinery (1975)
- [2] Schöniger, M., Waterman, M.S.: A local algorithm for DNA sequence alignment with inversions. *Bulletin of Mathematical Biology* **54** (1992) 521–536
- [3] Waterman, Eggert: A new algorithm for best subsequence alignments with application to tRNA-rRNA comparisons. *Journal of Molecular Biology* **197** (1987) 723–728
- [4] Kececioglu, J., Sankoff, D.: Exact and approximation algorithms for sorting by reversals, with application to genome rearrangement. *Algorithmica* **13** (1995) 180–210
- [5] Christie, D.A.: A $3/2$ -approximation algorithm for sorting by reversals. In: Proceedings of the Ninth Annual ACM-SIAM Symposium on Discrete Algorithms (San Francisco, CA, 1998), New York, ACM (1998) 244–252
- [6] Caprara, A.: Sorting permutations by reversals and Eulerian cycle decompositions. *SIAM J. Discrete Math.* **12** (1999) 91–110 (electronic)

- [7] Hannenhalli, S., Pevzner, P.A.: Transforming cabbage into turnip: polynomial algorithm for sorting signed permutations by reversals. In: ACM Symposium on Theory of Computing, Association for Computing Machinery (1995) 178–189
- [8] Hannenhalli, S., Pevzner, P.A.: Transforming cabbage into turnip: polynomial algorithm for sorting signed permutations by reversals. *J. ACM* **46** (1999) 1–27
- [9] Kaplan, H., Shamir, R., Tarjan, R.E.: Faster and simpler algorithm for sorting signed permutations by reversals. In: Proceedings of the Eighth Annual ACM-SIAM Symposium on Discrete Algorithms (New Orleans, LA, 1997), New York, ACM (1997) 344–351
- [10] Kaplan, H., Shamir, R., Tarjan, R.E.: A faster and simpler algorithm for sorting signed permutations by reversals. *SIAM J. Comput.* **29** (2000) 880–892 (electronic)
- [11] Bader, D.A., Moret, B.M.E., Yan, M.: A linear-time algorithm for computing inversion distance between signed permutations with an experimental study. In: Algorithms and data structures (Providence, RI, 2001). Volume 2125 of Lecture Notes in Comput. Sci. Springer, Berlin (2001) 365–376
- [12] El-Mabrouk, N.: Genome rearrangement by reversals and insertions/deletions of contiguous segments. In: Combinatorial pattern matching (Montreal, QC, 2000). Volume 1848 of Lecture Notes in Comput. Sci. Springer, Berlin (2000) 222–234
- [13] El-Mabrouk, N.: Sorting signed permutations by reversals and insertions/deletions of contiguous segments. *J. Discrete Algorithms* **1** (2000) 105–121
- [14] El-Mabrouk, N.: Reconstructing an ancestral genome using minimum segments duplications and reversals. *J. Comput. System Sci.* **65** (2002) 442–464 Special issue on computational biology 2002.
- [15] Hunt, J.W., Szymanski, T.G.: A fast algorithm for computing longest common subsequences. *Comm. ACM* **20** (1977) 350–353

- [16] Simon, I.: Sequence comparison: some theory and some practice. In Gross, M., Perrin, D., eds.: *Electronic Dictionaries and Automata in Computational Linguistics*, Berlin, Springer-Verlag (1989) 79–92 *Lecture Notes in Computer Science*, 377.
- [17] Eppstein, D., Galil, Z., Giancarlo, R., Italiano, G.F.: Sparse dynamic programming. I. Linear cost functions. *J. Assoc. Comput. Mach.* **39** (1992) 519–545
- [18] Eppstein, D., Galil, Z., Giancarlo, R., Italiano, G.F.: Sparse dynamic programming. II. Convex and concave cost functions. *J. Assoc. Comput. Mach.* **39** (1992) 546–567
- [19] Muchnik, I., do Lago, A.P., Llaca, V., Linton, E., Kulikowski, C.A., Messing, J.: Assignment-like optimization on bipartite graphs with ordered nodes as an approach to the analysis of comparative genomic data. *DIMACS Workshop on Whole Genome Comparison* (2001) <http://dimacs.rutgers.edu/Workshops/WholeGenome/>.
- [20] do Lago, A.P.: A sparse dynamic programming algorithm for alignment with inversions. *Workshop on Combinatorics, Algorithms, and Applications* (2003)