

A very elementary presentation of the Hannenhalli–Pevzner theory

Anne Bergeron

LACIM, Université du Québec à Montréal, C.P. 8888 Succ. Centre-Ville, Montréal, Que., Canada H3C 3P8

Received 28 February 2002; received in revised form 7 February 2003; accepted 23 April 2004

Available online 15 December 2004

Abstract

In 1995, Hannenhalli and Pevzner gave a first polynomial solution to the problem of finding the minimum number of reversals needed to sort a signed permutation. Their solution, as well as subsequent ones, relies on many intermediary constructions, such as simulations with permutations on $2n$ elements, and manipulation of various graphs. Here we give the first completely elementary treatment of this problem. We characterize *safe reversals* and *hurdles* working directly on the original signed permutation. Moreover, our presentation leads to polynomial algorithms that can be efficiently implemented using bit-wise operations. © 2004 Elsevier B.V. All rights reserved.

Keywords: Reversal distance; Signed permutations

1. Introduction

In the last 10 years, beginning with [7], many papers have been devoted to the subject of computing the *reversal distance* between two permutations. A *reversal* $\rho(i, j)$ transforms a permutation

$$\begin{aligned} \pi &= (\pi_1 \cdots \pi_i \pi_{i+1} \cdots \pi_j \cdots \pi_n) \\ \text{to } \pi' &= (\pi_1 \cdots \pi_j \cdots \pi_{i+1} \pi_i \cdots \pi_n) \end{aligned}$$

and the *reversal distance* between two permutations is the minimum number of reversals that transform one into the other.

From a problem of unknown complexity, it graduated to an NP-Hard problem [3], but an interesting variant was proven to be polynomial [4]. In the *signed* version of the problem, each element of the permutation has a plus or minus sign, and a reversal $\rho(i, j)$ transforms π to

$$\pi' = (\pi_1 \cdots -\pi_j \cdots -\pi_{i+1} -\pi_i \cdots \pi_n).$$

Permutations, and their reversals, are useful tools in the comparative study of genomes. The genome of a species can be thought of as a set of ordered sequences of genes, the ordering devices being the chromosomes, each gene having an orientation given by its location on the DNA double strand. Different species often share similar genes that were inherited from common ancestors. However, these genes have been shuffled by mutations that modified the content of chromosomes, the order of genes within a particular chromosome, and/or the orientation of a gene. Comparing two sets of similar genes appearing along a chromosome in two different species yields two (signed) permutations. It is widely accepted that the reversal distance between these two permutations provides a good estimate of the evolutionary distance between the two species.

E-mail address: bergeron.anne@uqam.ca (A. Bergeron).

Computing the reversal distance between signed permutations is a delicate task, since some reversals unexpectedly affect deep structures in permutations. In 1995, Hannenhalli and Pevzner proposed the first polynomial algorithm to solve it [4], developing along the way a theory of how and why some permutations were particularly resistant to sorting by reversals. It is of no surprise that the label *fortress* was assigned to specially acute cases.

Hannenhalli and Pevzner relied on several intermediate constructions that have been simplified since [1,2,5], but grasping all the details remains a challenge. All the criteria given for choosing a *safe* reversal involve the construction of an associate permutation on $2n$ points, and the analysis of cycles and/or connected component of graphs associated with this permutation.

In this paper, we present an elementary treatment of the sorting of the *oriented components* of a permutation, together with a new definition of the concept of *hurdle*, that further simplifies the definition given in [5]. Our first algorithm is so simple that, for example, sorting a permutation of length 20, *by hand*, should be easy and straightforward.

The next section presents the basic algorithms. Section 3 contains the necessary links to the Hannenhalli–Pevzner theory, and the proofs of the claims in Section 2. Finally, in the last section, we discuss complexity issues, and we give a *bit-vector* implementation of the sorting algorithm that runs in $\mathcal{O}(n^2)$ bit-vector operations, or in $\mathcal{O}(n^3/w)$ operations, where w is the word-size of the processor.

2. Basic sorting

The problem of computing the distance between two permutations is often recast as the problem of computing $d(\pi)$, the reversal distance between a permutation π and the identity permutation $(1\ 2\ \dots\ n)$. In this paper, we focus on the reconstruction of one possible sequence of reversals that realizes $d(\pi)$, also called the *sorting by reversals problem*. As usual, we will assume that the permutation is *framed* by 0 and $n + 1$, and that those extra elements are always positive:

$$\pi = (0\ \pi_1\ \pi_2\ \dots\ \pi_n\ n + 1).$$

An *oriented pair* (π_i, π_j) , $i < j$, is a pair of consecutive integers, that is $|\pi_i| - |\pi_j| = \pm 1$, with opposite signs. For example, the oriented pairs of the permutation

$$(0\ 3\ 1\ 6\ 5\ -2\ 4\ 7)$$

are $(1, -2)$ and $(3, -2)$.

Oriented pairs are useful, in the sense that they indicate reversals that create consecutive elements of the permutation that are also consecutive integers. For example, the pair $(1, -2)$ induces the reversal

$$(0\ 3\ 1\ \underline{6\ 5}\ -2\ 4\ 7),$$

$$(0\ 3\ 1\ 2\ -5\ -6\ 4\ 7),$$

creating the consecutive integers 1 2. Note that the reverse of a pair of consecutive integers, such as $-2\ -1$, is also a pair of consecutive integers.

In general, the reversal induced by an oriented pair (π_i, π_j) will be

$$\rho(i, j - 1), \quad \text{if } \pi_i + \pi_j = +1 \quad \text{and}$$

$$\rho(i + 1, j), \quad \text{if } \pi_i + \pi_j = -1.$$

Note that reversals that create consecutive integers are always induced by oriented pairs. Such a reversal is called an *oriented reversal*. We define the *score* of an (oriented) reversal as the number of oriented pairs in the resulting permutation. For example, the score of the reversal

$$(0\ \underline{3\ 1\ 6\ 5}\ -2\ 4\ 7),$$

$$(0\ -5\ -6\ -1\ -3\ -2\ 4\ 7)$$

is 4, since the resulting permutation has four oriented pairs. Computing the score of a reversal is tedious but elementary, and we will discuss efficient algorithms too do so in Section 4. The fact that oriented reversals have a beneficial effect on the ordering of a permutation suggests a first sorting strategy.

Algorithm 1. *As long as π has an oriented pair, choose the oriented reversal that has maximal score.*

For example, the scores of two oriented pairs $(1, -2)$ and $(3, -2)$ of the permutation

$$(0 \quad \underline{3 \quad 1 \quad 6 \quad 5} \quad -2 \quad 4 \quad 7)$$

are, respectively, 2 and 4. So we choose the reversal induced by $(3, -2)$, yielding the new permutation

$$(0 \quad -5 \quad -6 \quad -1 \quad \underline{-3 \quad -2} \quad 4 \quad 7).$$

This permutation has now four oriented pairs $(0, -1)$, $(-3, 4)$, $(-5, 4)$ and $(-6, 7)$, all of which have score 2, except $(-3, 4)$. Acting on this pair yields

$$(0 \quad -5 \quad -6 \quad \underline{-1} \quad 2 \quad 3 \quad 4 \quad 7)$$

which has four oriented pairs. Note here that the score of the pair $(0, -1)$ is 0. The corresponding oriented reversal would produce a permutation with no oriented pair, and the algorithm would stop, in this case with an unsorted permutation. Fortunately, the pair $(-1, 2)$ has a positive, and maximal, score and we get, in a similar way, the last two necessary reversals to sort the permutation

$$(0 \quad -5 \quad \underline{-6 \quad 1 \quad 2 \quad 3 \quad 4} \quad 7),$$

$$(0 \quad \underline{-5 \quad -4 \quad -3 \quad -2 \quad -1} \quad 6 \quad 7),$$

$$(0 \quad 1 \quad 2 \quad 3 \quad 4 \quad 5 \quad 6 \quad 7).$$

Interestingly enough, this elementary strategy is sufficient to optimally sort most random permutations and almost all permutations that arise from biological data. The strategy is also optimal, and we will prove in the next section the following claim.

Claim 1. *If the strategy of Algorithm 1 applies k reversals to a permutation π , yielding a permutation π' , then $d(\pi) = d(\pi') + k$.*

The output of Algorithm 1 will be a permutation of positive elements. Most reversal applied to such permutations will create oriented pairs, but the choice of an optimal reversal is delicate. We discuss this problem in the next paragraph.

2.1. Sorting positive permutations

Let π be a signed permutation with only positive elements, and assume that π is *reduced*, that is π does not contain consecutive integers. Suppose also that π is framed by 0 and $n + 1$ and consider, as in [5], the circular order induced by setting 0 to be the successor of $n + 1$.

Define a *framed interval* in π as an interval of the form

$$i \quad \pi_{j+1} \quad \pi_{j+2} \quad \cdots \quad \pi_{j+k-1} \quad i + k,$$

such that all integers between i and $i + k$ belong to the interval $[i \cdots i + k]$. For example, consider the permutation

$$(0 \quad 2 \quad 5 \quad 4 \quad 3 \quad 6 \quad 1 \quad 7).$$

The whole permutation is a framed interval. But we have also the interval: $2 \ 5 \ 4 \ 3 \ 6$, which can be reordered as $2 \ 3 \ 4 \ 5 \ 6$, and, by circularity, the interval $6 \ 1 \ 7 \ 0 \ 2$, which can be reordered as $6 \ 7 \ 0 \ 1 \ 2$, since 0 is the successor of 7.

Definition 2. If π is reduced, a *hurdle* in π is a framed interval that contains no shorter framed interval.

Claim 3. *Hurdles as defined in Definition 2 are the same hurdles that are defined in [4,5].*

When a permutation has only one or two hurdles, 1 reversal is sufficient to create enough oriented pairs in order to completely sort the permutation with Algorithm 1. Two operations are introduced in [4], the first one is *hurdle cutting* which consist in reversing the segment between i and $i + 1$ of a hurdle

$$\underline{i \quad \pi_{j+1} \quad \pi_{j+2} \quad \cdots \quad i + 1 \quad \cdots \quad \pi_{j+k-1} \quad i + k}.$$

This reversal is sufficient to sort all the interval using Algorithm 1. For example, the following permutation contains only one hurdle

$$(0 \quad 2 \quad 4 \quad 3 \quad 1 \quad 5).$$

The reversal of elements 2, 4 and 3 cuts the hurdle, and the resulting permutation

$$(0 \quad -3 \quad -4 \quad -2 \quad 1 \quad 5)$$

can be sorted with 4 reversals by Algorithm 1.

The second operation is *hurdle merging*, which acts on the endpoints of two hurdles:

$$i \cdots \underline{i+k} \cdots i' \cdots i'+k',$$

and does the reversal $\rho(i+k, i')$. If a permutation has only two hurdles, merging them will produce a permutation that can be completely sorted by Algorithm 1.

Thus, for example, merging the two hurdles in the permutation

$$(0 \quad 2 \quad 5 \quad 4 \quad 3 \quad 6 \quad 1 \quad 7)$$

yields the permutation

$$(0 \quad 2 \quad 5 \quad 4 \quad 3 \quad -6 \quad 1 \quad 7),$$

which can be sorted in 5 reversals using Algorithm 1.

Merging and cutting hurdles in a permutation that contains more than two hurdles must be managed carefully. Indeed, cutting some hurdles can create new ones!

Definition 4. A simple hurdle is a hurdle whose cutting decreases the number of hurdles. Hurdles that are not simple are called *super hurdles*.

For example, the permutation $(0 \ 2 \ 5 \ 4 \ 3 \ 6 \ 1 \ 7)$ has two hurdles. Cutting and sorting the hurdle $2 \ 5 \ 4 \ 3 \ 6$ yields the permutation,

$$(0 \quad 2 \quad 3 \quad 4 \quad 5 \quad 6 \quad 1 \quad 7),$$

which, by collapsing the sequence $2 \ 3 \ 4 \ 5 \ 6$, reduces to

$$(0 \quad 2 \quad 1 \quad 3),$$

which has only one hurdle.

However, the permutation $(0 \ 2 \ 4 \ 3 \ 5 \ 1 \ 6 \ 8 \ 7 \ 9)$ also contains two hurdles, and when one cuts the hurdle $2 \ 4 \ 3 \ 5$, the resulting reduced permutation is

$$(0 \quad 2 \quad 1 \quad 3 \quad 5 \quad 4 \quad 6),$$

which still has two hurdles.

The following algorithm is adapted from [5], and is discussed originally in [4].

Algorithm 2. *If a permutation has $2k$ hurdles, $k \geq 2$, merge any two non-consecutive hurdles. If a permutation has $2k + 1$ hurdles, $k \geq 1$, then if it has one simple hurdle, cut it; If it has none, merge two non-consecutive hurdles, or consecutive ones if $k = 1$.*

Together with Algorithm 1, Algorithm 2 can be used to optimally sort any signed permutation. Permutations that are not already reduced can always be reduced by merging consecutive integers, and by renumbering all the elements.

This completes the first part of the paper, and, in the next section, we turn to the task of proving our various claims.

3. Selected results from the Hannenhalli–Pevzner theory

The exposition of the complete results of the Hannenhalli–Pevzner theory is beyond the scope of this paper, and the reader is referred to the original paper [4], or the book on computational molecular biology by Pevzner [6]. Instead, we will show the soundness of our algorithms by directly using the *arc overlap graph* introduced in [5].

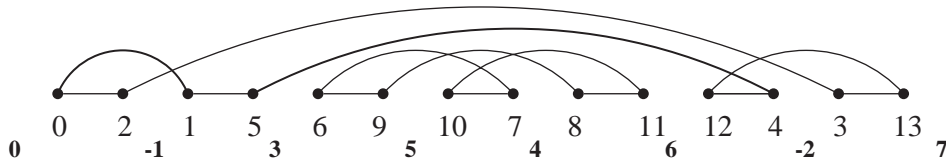


Fig. 1. The Breakpoint Graph of $\pi = (0 -1 3 5 4 6 -2 7)$.

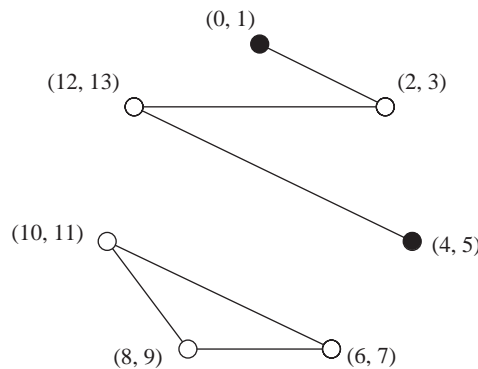


Fig. 2. The arc overlap graph of $\pi = (0 -1 3 5 4 6 -2 7)$.

The first construction is the *breakpoint graph* associated with π . Each positive element x in the permutation π is replaced by the sequence $2x - 1 \ 2x$, and each negative element $-x$ by the sequence $2x \ 2x - 1$. For example,

$$\pi = (0 \ -1 \ 3 \ 5 \ 4 \ 6 \ -2 \ 7)$$

becomes

$$\pi' = (0 \ 2 \ 1 \ 5 \ 6 \ 9 \ 10 \ 7 \ 8 \ 11 \ 12 \ 4 \ 3 \ 13).$$

Reversals $\rho(i, j)$ of π are simulated by unsigned reversals $\rho(2i - 1, 2j)$ in π' .

The elements of π' are the vertices of the breakpoint graph. Straight edges join every other pair of consecutive elements of π' , starting with 0, and curved edges, called *arcs*, join every other pair of consecutive integers, starting with $(0, 1)$.

Every connected component of the breakpoint graph is a cycle, which is a consequence of the fact that each vertex has exactly two incident edges. The graph of Fig. 1 has 2 cycles.

The *support* of an arc is the interval of elements of π' between, and including, its endpoints. Two arcs *overlap* if their support intersect, without proper containment. An arc is *oriented* if its support contains an odd number of elements, otherwise it is *unoriented*. In Fig. 1, the oriented arcs are $(0, 1)$ and $(4, 5)$. Note that an arc is oriented if and only if its endpoints are images of an oriented pair of the original permutation.

The *arc overlap graph* is the graph whose vertices are arcs of the breakpoint graph, and whose edges join overlapping arcs. The overlap graph corresponding to the breakpoint graph of Fig. 1 is illustrated in Fig. 2, in which each vertex is labeled by an arc $(2x, 2x + 1)$. *Oriented* vertices—those for which the corresponding arc is oriented—are marked by black dots. Orientation extends to connected components in the sense that a connected component with at least one oriented vertex is oriented. It is easy to show that a vertex is oriented if and only if its degree is odd (see Lemma 5).

There is a natural bijection between the vertices of the overlap graph and pairs of (unsigned) consecutive integers x and $x + 1$ in the original permutation. Indeed, a pair of consecutive integers will generate four consecutive integers in the unsigned permutation: $2x - 1, 2x, 2x + 1,$ and $2x + 2$. The vertex $(2x, 2x + 1)$ is associated with the pair x and $x + 1$. For example, the oriented pair $(3, -2)$ in π corresponds to the vertex $(4, 5)$ in the overlap graph. We will refer to the *reversal induced by a vertex* meaning the reversal induced by the corresponding oriented pair of the original permutation. The following lemmas, mostly from [5], pinpoint the important relations between a signed permutation and its overlap graph, and will help to prove Claim 1 of Section 2.

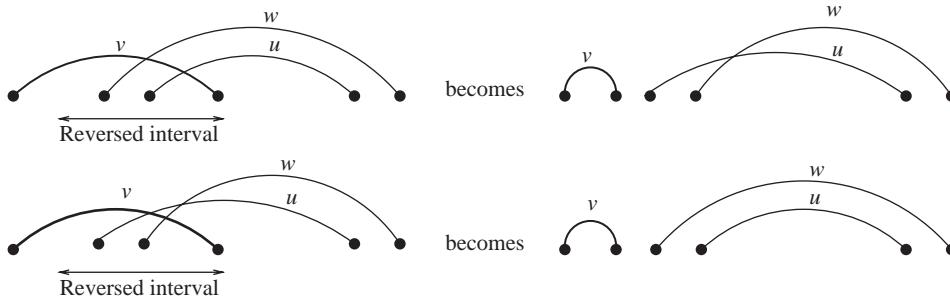


Fig. 3. Complementations of subgraphs.

Lemma 5. *A vertex has an odd degree if and only if it is oriented.*

Proof. Let $2x - 1, 2x, 2x + 1,$ and $2x + 2,$ be the four integers associated with the oriented pair (π_i, π_j) . Since π_i and π_j have different signs, the positions of $2x$ and $2x + 1$ will not have the same parity in the unsigned permutation. Thus, the interval between $2x$ and $2x + 1$ has an odd length, implying that it overlaps an odd number of other intervals. On the other hand, any interval that overlaps an odd number of intervals must have an odd length. Therefore, the positions of its endpoints must have different parities, implying that the corresponding pair of consecutive integers is oriented. \square

Lemma 6. *If one performs the reversal corresponding to an oriented vertex v , the effect on the overlap graph will be to complement the subgraph induced by v and its adjacent vertices.*

Proof. The reversal corresponding to an oriented vertex v has the effect of collapsing the associated interval, thus v will become isolated. Let u and w be two intervals overlapping v , meaning that exactly one of their endpoints lies in the interval spanned by v . The reversal induced by v will reverse these two points. Here, a picture (Fig. 3) is worth a thousand words. \square

Lemma 7. *If one performs the reversal corresponding to an oriented vertex v , each vertex adjacent to v will change its orientation.*

Proof. Since v is oriented, it has an odd number $2k + 1$ of adjacent vertices. Let w be a vertex adjacent to v , with j neighbors also adjacent to v . With the reversal, w will lose $j + 1$ neighbors, and gain $2k - j$ new ones. Thus the degree of w will change by $2k - 2j - 1$, changing its orientation. \square

Lemma 8. *The score of the oriented reversal corresponding to an oriented vertex v is given by*

$$T + U - O - 1,$$

where T is the total number of oriented vertices in the graph, U is the number of unoriented vertices adjacent to v , and O is the number of oriented vertices adjacent to v .

Proof. This follows trivially from the preceding lemmas. \square

We now state a basic result that is proved, in different ways, in [4,5]. Define an *oriented component* of the overlap graph as a connected component that contains at least one oriented vertex, otherwise the component is *unoriented*. A *safe reversal* is a reversal that does not create new unoriented components, except for isolated vertices. The following theorem states that safe reversals, when they exist, always decrease the distance by 1.

Theorem 9 (Hannenhalli and Pevzner [4]). *Any sequence of oriented safe reversals is optimal.*

The difficulties in sorting oriented components lie in the detection of safe reversals. Hannenhalli and Pevzner deal with the problem by computing several statistics on cycles and breakpoints of various graphs. Kaplan et al. [5] solve it by searching for particular cliques in the overlap graph. The next theorem argues that the elementary strategy of choosing the reversal with maximal score is optimal, thus proving Claim 1.

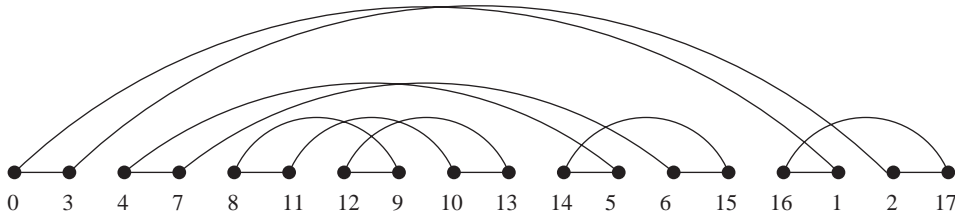


Fig. 4. The breakpoint graph of $\pi = (0\ 2\ 4\ 6\ 5\ 7\ 3\ 8\ 1\ 9)$.

Theorem 10. *An oriented reversal of maximal score is safe.*

Proof. Suppose that vertex v has maximal score, and that the reversal induced by v creates a new unoriented component C containing more than one vertex. At least one of the vertices in C must have been adjacent to v , since the only edges affected by the reversal are those between vertices adjacent to v . Let w be a vertex formerly adjacent to v and contained in C , and consider the scores of v and w :

$$score(v) = T + U - O - 1,$$

$$score(w) = T + U' - O' - 1.$$

All unoriented vertices formerly adjacent to v must have been adjacent to w . Indeed, an unoriented vertex adjacent to v and not to w will become oriented, and connected to w , contrary to the assumption that C is unoriented. Thus, $U' \geq U$.

All oriented vertices formerly adjacent to w must have been adjacent to v . If this was not the case, an oriented vertex adjacent to w but not to v would remain oriented, again contradicting the fact that C is unoriented. Thus, $O' \leq O$.

Now, if both $O' = O$ and $U' = U$, vertices v and w have the same set of adjacent vertices, and complementing the subgraph of v and its adjacent vertices will isolate both v and w . Therefore, we must have that $score(w) > score(v)$, which is a contradiction. \square

3.1. Hurdles

The goal of this section is to prove Claim 3. We assume that π is a positive and reduced permutation. These assumptions are equivalent to saying that the overlap graph has no oriented components, all of which can be cleared by Algorithm 1, and no isolated vertices.

Consider again the circular order, this time on the interval $[0..2n - 1]$, induced by setting 0 to be the successor of $2n - 1$. The span of a set of vertices X in the overlap graph is the minimum interval that contains, in the circular order, all vertices in X . For example, consider the permutation

$$\pi = (0\ 2\ 4\ 6\ 5\ 7\ 3\ 8\ 1\ 9),$$

whose breakpoint graph and arc overlap graphs are illustrated in Figs. 4 and 5. The three connected components of its arc overlap graph have spans:

$$[4, 15] = [4\ 7\ 8\ 11\ 12\ 9\ 10\ 13\ 14\ 5\ 6\ 15]$$

$$[8, 13] = [8\ 11\ 12\ 9\ 10\ 13]$$

$$[16, 3] = [16\ 1\ 2\ 17\ 0\ 3].$$

Hurdles are defined in [4] as unoriented components which are minimal with respect to the order induced by span inclusion. Moreover, in [5], it is shown that the span of a connected component is always of the form $[2i, 2j - 1]$. The following lemmas and theorem detail the relationships between connected components and framed intervals, substantiating the second claim of Section 2.

Lemma 11. *Framed intervals of the form $[i, j]$ in a permutation on n elements are in one-to-one correspondence with framed intervals of the form $[2i, 2j - 1]$ in the corresponding unsigned permutation on $2n$ elements.*

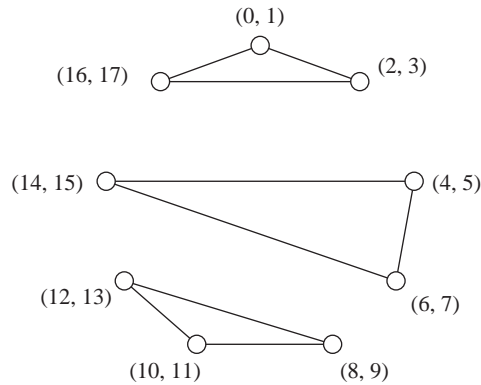


Fig. 5. The arc overlap graph of $\pi = (0\ 2\ 4\ 6\ 5\ 7\ 3\ 8\ 1\ 9)$.

Proof. The endpoints i and j of a framed interval $[i, j]$ will be mapped, respectively, to the pairs $2i - 1, 2i$, and $2j - 1, 2j$. All integers between i and j appear in the interval $[i, j]$, if and only if all the integers between $2i$ and $2j - 1$ appear in the interval $[2i, 2j - 1]$. \square

Lemma 12. Any framed interval $[2i, 2j - 1]$ is the span of a union of connected components.

Proof. If $[2i, 2j - 1]$ is a framed interval, it contains exactly the integers between $2i$ and $2j - 1$, thus the only arcs in this interval are: $(2i, 2i + 1), (2i + 2, 2i + 3), \dots, (2j - 2, 2j - 1)$, and no other arc intersects this set. Therefore, the corresponding set of vertices is not connected to any other vertex. \square

Lemma 13. The span $[2i, 2j - 1]$ of a connected component is always a framed interval.

Proof. If vertex $(2i, 2i + 1)$ is connected to $(2j - 2, 2j - 1)$, there must be a sequence of intersecting arcs linking $2i$ to $2j - 1$.



Any arc with only one endpoint between $2i$ and $2j - 1$ would therefore intersect one of the arcs in the sequence, and would be part of the connected component, so there are none. Thus, if integer $2k$ is in the interval, then $2k + 1$ is also in the interval, and if $i \leq k < j - 1$, then $2k + 2$ is also in the interval. \square

Theorem 14. If π is reduced, an unoriented component is minimal if and only if its span is a framed interval that contains no other.

Proof. By Lemma 13, the span of a connected component is always a framed interval. If the component is minimal with respect to span inclusion, by Lemma 12, its span cannot contain properly another framed interval.

On the other hand, a framed interval $[2i, 2j - 1]$ that contains no other yields a single connected component C whose vertices endpoints are exactly the integers between $2i$ and $2j - 1$. Thus the vertices of C are consecutive on the circle, and component C is minimal. \square

Using Lemma 11, Theorem 14 gives an elementary characterization of the concept of hurdles, that does not rely on the construction of the overlap graph.

4. Settling scores

Algorithm 1 has a straightforward naive implementation: in order to find a safe reversal, perform each possible oriented reversal on the original permutation, and count the number of resulting oriented reversals. Since there are $\mathcal{O}(n)$ oriented reversal

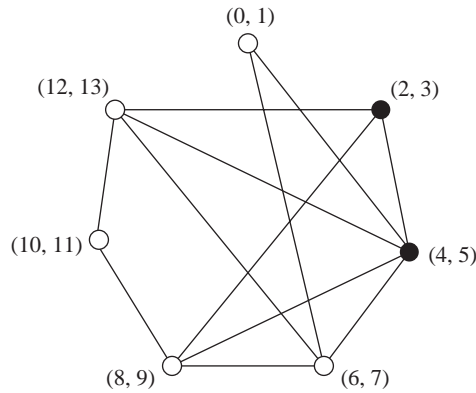


Fig. 6. The arc overlap graph of (0 3 1 6 5 - 2 4 7).

in a given permutation on n elements, and since computing the number of oriented reversals in a permutation can be done in $\mathcal{O}(n)$ operations, this yields an $\mathcal{O}(n^3)$ algorithm for sorting a permutation on n elements.

However, performing a reversal, or computing its score, is a “local” operation on the overlap graph, and this locality suggests the possibility of a parallel algorithm to keep the scores and to compute the effects of a reversal. The parallelism of the algorithm will exploit the inherent parallelism of basic operations of a processor. We will work with bit-vectors, and use only three operations on these vectors: the exclusive-or operator \oplus ; the conjunction \wedge ; and the negation $-$.

Subsets of vertices of the arc overlap graph will be represented by *characteristics* bit-vectors: If s is a subset of the $\{0, \dots, n\}$, then the bold symbol s is the bit-vector:

$$s = (s_0, \dots, s_n) \quad \text{where } s_i = \begin{cases} 1 & \text{if } i \in s, \\ 0 & \text{otherwise.} \end{cases}$$

4.1. The data structure

Given an arc overlap graph, we first construct a bit-matrix in which each line v_i is the set of adjacent vertices to arc $(2i, 2i + 1)$. For example, consider the permutation (0 3 1 6 5 - 2 4 7), whose arc overlap graph is illustrated in Fig. 6.

The bit-matrix associated with this graph is the following:

	v_0	v_1	v_2	v_3	v_4	v_5	v_6
v_0	0	0	1	1	0	0	0
v_1	0	0	1	0	1	0	1
v_2	1	1	0	1	1	0	1
v_3	1	0	1	0	1	0	1
v_4	0	1	1	1	0	1	0
v_5	0	0	0	0	1	0	1
v_6	0	1	1	1	0	1	0
p	0	1	1	0	0	0	0
s	0	1	3	2	0	2	0

The last two lines contain, respectively, the parity p , or orientation, of the vertex, and the score $s = U - O$ of the associated reversal. [The score of a reversal was defined in Section 3 as $T + U - O - 1$, where T is the total number of oriented reversals in the original permutation. Since we want to maximize scores, and since $T - 1$ is constant for all reversals, it is only necessary to maximize $U - O$.]

We will discuss efficient ways to initialize the structure and to adjust scores in Sections 4.2 and 4.3.

Given the vectors \mathbf{p} and \mathbf{s} , selecting the oriented reversal with maximal score is elementary. In the above example, vertex 2 would be the selected candidate.

The interesting part is how a reversal affects the structure. These effects are summarized in the following algorithm, which recalculates the bit-matrix \mathbf{v} , the parity vector \mathbf{p} , and the score vector \mathbf{s} , following the reversal associated to vertex i , whose set of adjacent vertices is denoted by \mathbf{v}_i .

1. $\mathbf{s} \leftarrow \mathbf{s} + \mathbf{v}_i$
2. $\mathbf{v}_{ii} \leftarrow 1$
3. **for** each vertex j adjacent to i **do**
4. **if** j is oriented
5. $\mathbf{s} \leftarrow \mathbf{s} + \mathbf{v}_j$
6. $\mathbf{v}_{jj} \leftarrow 1$
7. $\mathbf{v}_j \leftarrow \mathbf{v}_j \oplus \mathbf{v}_i$
8. $\mathbf{s} \leftarrow \mathbf{s} + \mathbf{v}_j$
9. **else**
10. $\mathbf{s} \leftarrow \mathbf{s} - \mathbf{v}_j$
11. $\mathbf{v}_{jj} \leftarrow 1$
12. $\mathbf{v}_j \leftarrow \mathbf{v}_j \oplus \mathbf{v}_i$
13. $\mathbf{s} \leftarrow \mathbf{s} - \mathbf{v}_j$
14. $\mathbf{p} \leftarrow \mathbf{p} \oplus \mathbf{v}_i$

The logic behind the algorithm is the following. Since vertex i will become unoriented and isolated, each vertex adjacent to i will automatically gain a point of score in step 1. Next, if j is a vertex adjacent to i , vertices adjacent to j after the reversal are either existing vertices that were adjacent to j and not adjacent to i , or vertices that were adjacent to i but not to j . This is the definition of the exclusive-or operator \oplus .

The exceptions to this rule are i and j themselves, and this problem is solved by setting the diagonal bits to 1 before computing $\mathbf{v}_j \oplus \mathbf{v}_i$ in steps 6 and 11.

If j is oriented, each of its former adjacent vertices will gain one point of score, since j will become unoriented, and each of its new adjacent vertices will gain one point of score. We thus add these points of score with two instructions sandwiching step 6. Note that a vertex that *stays* connected to j will gain a total of two points. For unoriented vertices, the gains are converted to losses.

The amount of work done to process a reversal corresponding to vertex i , in terms of vector operations, is thus proportional to the number of adjacent vertices to vertex i .

4.2. Representing the scores

The additions and subtractions to adjust the score vector are the usual arithmetic operations performed component-wise. In order to have a truly bit-vector implementation, we represented the score vector \mathbf{s} as a $\lceil \log(n) \rceil \times n$ bit-matrix, each column containing the binary representation of a score. The k th line of the matrix is referred to as s_k . With this representation, component-wise addition of a bit-vector \mathbf{v} to \mathbf{s} can be realized with the following:

1. **for** k **from** 1 **to** $\lceil \log(n) \rceil$
2. $\mathbf{t} \leftarrow \mathbf{v}$
3. $\mathbf{v} \leftarrow \mathbf{v} \wedge s_k$
4. $s_k \leftarrow \mathbf{t} \oplus s_k$

Subtraction is implemented in a similar way. A side benefit of this structure is that the selection of the next reversal can be also done in parallel, by “sifting” the score matrix through the parity vector. The set \mathbf{c} of candidates contains initially all the oriented vertices. Going from the higher bit of scores to the lower, if at least one of the candidates has bit i set to 1, we eliminate all candidates for which bit i is 0.

1. $\mathbf{c} \leftarrow \mathbf{p}$
2. $i \leftarrow \lceil \log(n) \rceil$
3. **while** $i \geq 0$ **do**
4. **while** $(\mathbf{c} \wedge s_i) = 0$
5. $i \leftarrow i - 1$
6. **if** $i \geq 0$
7. $\mathbf{c} \leftarrow \mathbf{c} \wedge s_i$
8. $i \leftarrow i - 1$

At the end of the loop, c is the set of oriented vertices of maximal score.

4.3. Initializing the data structure

We saw, in Section 3, that the overlap graph of a signed permutation $\pi = (\pi_1 \pi_2 \cdots \pi_n)$ contains $n + 1$ vertices corresponding to the arcs joining $2x$ and $2x + 1$ in the equivalent unsigned permutation. In this section, we will construct a representation of the overlap graph without explicitly referring to the unsigned permutation, thus removing one more step between the actual algorithm, and the original formulation of the problem.

The construction is based on the following simple lemma. Let I be a set of intervals with distinct endpoints in an ordered set S . Denote an interval $i \in I$ by the ordered pair (b_i, e_i) of its endpoints. Define the sets l_i and r_i as follows:

$$r_i = \{j \in I \mid b_j < e_i < e_j\},$$

$$l_i = \{j \in I \mid b_j < b_i < e_j\}.$$

The set r_i is the set of intervals j in I that contains the right endpoint of interval i , and the set l_i is the set of intervals j in I that contains the left endpoint of interval i . We have the following.

Lemma 15. *The set v_i of intervals that overlap i in I is given by: $v_i = l_i \oplus r_i$.*

Starting with a signed permutation $\pi = (\pi_1 \pi_2 \cdots \pi_n)$, we first read the elements from left to right. Let a represent the set of arcs for which exactly one endpoint has been read. Initially, a is the set $\{0\}$, corresponding to the arc $(0, 1)$. When element π_i is read, we have to process two arcs: $(2\pi_i - 2, 2\pi_i - 1)$ and $(2\pi_i, 2\pi_i + 1)$. In increasing order, if π_i is positive, and decreasing order, otherwise. Processing an arc $(2j, 2j + 1)$ is done by the following instructions:

1. **if** $a_j = 0$
2. $a_j \leftarrow 1$ (* First endpoint of arc $(2j, 2j + 1)$ *)
3. **else** (* Second endpoint *)
4. $a_j \leftarrow 0$
5. $v_j \leftarrow a$ (* a is the set r_j *)

We then repeat the process in the reverse order, reading the permutation from right to left, initializing a to the set $\{n\}$, and changing the last instruction to $v_j \leftarrow v_j \oplus a$. The parity vector p is initialized in this loop with the instruction $p \leftarrow p \oplus v_j$.

Once the bit-matrix v is computed, the scores are initialized by either adding or subtracting each v_j , depending on its parity, to an initial null matrix of scores.

4.4. Analysis

The formal analysis of the algorithm of Section 4.1 raises interesting questions. For example, what is an elementary operation? Except for a few control statements, the only operations used by the algorithm are very efficient bit-wise logical operators on words of size w , typically 32 or 64, depending on implementation. The most expensive instructions in the main loop are additions and subtractions, such as

$$s \leftarrow s + v_j,$$

where s is a bit-matrix of size $n \log(n)$, and v_j is a bit-vector of size n . Such an operation requires a total of $(2n \log(n))/w$ elementary operations with the loop described in Section 4.2. Hopefully, $\log(n)$ is much smaller than w , and, in the range of biologically meaningful values, n is often a small multiple of w . In the actual implementation, the loop is controlled by the value of $\log(\text{maximal score})$ which tends to be much less than $\log(n)$. We thus have a, very generous, $\mathcal{O}(n)$ estimate for the instructions in the main loop.

The overall work done by the algorithm depends on the total number v of vertices adjacent to vertices of maximal score. We can easily bound it by n^2 , noting that the number d of reversals needed to sort the permutation is bounded by n , and the degree of a vertex is also bounded by n . We thus get an $\mathcal{O}(n^3)$ estimate for the algorithm, assuming that $\log n < w$. However, in practical applications, such as [8], this algorithm outperforms consistently classical algorithms with smaller theoretical complexity.

References

- [1] D. Bader, B. Moret, M. Yan, A linear-time algorithm for computing inversion distance between signed permutations with an experimental study, *J. Comput. Biol.* 8 (5) (2001) 483–491.

- [2] P. Berman, S. Hannenhalli, Fast Sorting by Reversal, CPM 1996, Lecture Notes in Comput. Sci., vol. 1075, Springer, Berlin, 1996, pp. 168–185.
- [3] A. Caprara, Sorting by Reversals is Difficult, RECOMB 1997, ACM Press, New York, 1997, pp. 75–83.
- [4] S. Hannenhalli, P.A. Pevzner, Transforming cabbage into turnip: polynomial algorithm for sorting signed permutations by reversals, *J. Assoc. Comput. Mach.* 46 (1) (1999) 1–27.
- [5] H. Kaplan, R. Shamir, R. Tarjan, A faster and simpler algorithm for sorting signed permutations by reversals, *SIAM J. Comput.* 29 (3) (1999) 880–892.
- [6] P. Pevzner, *Computational Molecular Biology*, MIT Press, Cambridge, MA, 2000, 314p.
- [7] D. Sankoff, Edit Distances for Genome Comparisons Based on Non-Local Operations, CPM 1992, Lecture Notes in Comput. Sci., vol. 644, Springer, Berlin, 1992, pp. 121–135.
- [8] A. Siepel, An Algorithm to Find all Sorting Reversals, RECOMB 2002, ACM Press, New York, 2002, pp. 281–290.