

# Faster and Simpler Algorithm for Sorting Signed Permutations by Reversals

Haim Kaplan\*

Ron Shamir†

Robert E. Tarjan‡

## Abstract

We give a quadratic algorithm for finding the minimum number of reversals needed to sort a signed permutation. Our algorithm is faster than the previous algorithm of Hannenhalli and Pevzner and its faster implementation of Berman and Hannenhalli. The algorithm is conceptually simple and does not require special data structures. Our study also considerably simplifies the combinatorial structures used by the analysis.

## 1 Introduction

In this paper we study the problem of sorting signed permutations by reversals. A *signed permutation* is a permutation  $\pi = (\pi_1, \dots, \pi_n)$  on the integers  $\{1, \dots, n\}$ , where each number is also assigned a sign of plus or minus. A *reversal*,  $\rho(i, j)$ , on  $\pi$  transforms  $\pi$  to

$$\pi' = \pi\rho(i, j) = (\pi_1, \dots, \pi_{i-1}, -\pi_j, -\pi_{j-1}, \dots, -\pi_i, \pi_{j+1}, \dots, \pi_n).$$

The minimum number of reversals needed to transform one permutation to another is called the *reversal distance* between them. The problem of *sorting signed permutations by reversals* is to find for a given signed permutation  $\pi$  its reversal distance from the identity permutation  $(+1, +2, \dots, +n)$ .

The motivation to studying the problem arises in molecular biology: Concurrent with the fast progress of the Human Genome Project, genetic and DNA data on many model organisms is accumulating rapidly,

and consequently the ability to compare genomes of different species has grown dramatically. One of the best ways of checking similarity between genomes on a large scale is to compare the order of appearance of identical genes in the two species. In the Thirties, Dobzhansky and Sturtevant [6] have already studied the notion of inversions in chromosomes of *drosophila*. Beginning in the late Eighties, Jeffrey Palmer has demonstrated that different species may have essentially the same genes, but the gene order may differ between species. Taking a distant perspective, the genes along a chromosome can be thought of as points along a line. Numbers identify the particular genes, and as genes have directionality, signs correspond to their direction. Palmer and others have shown that the difference in order may be explained by a small number of reversals [15, 16, 17, 18, 11]. These reversals correspond to evolutionary changes along the history between the two genomes, so the number of reversals reflects the evolutionary distance between the species. Hence, given two such permutations, their reversal distance measures their evolutionary distance.

Mathematical analysis of genome rearrangement problems was initiated by Sankoff [20, 19]. Kececioglu and Sankoff [14] gave the first constant factor polynomial approximation algorithm for the problem and conjectured that the problem is NP-hard. Bafna and Pevzner [2] have subsequently improved the approximation factor, and additional studies have revealed a rich combinatorial structure for rearrangement problems [13, 12, 3, 8, 10]. Quite recently, Caprara [5] has established that sorting *unsigned* permutations is indeed NP-hard, using some of the combinatorial tools developed by Bafna and Pevzner [2].

In 1995, Hannenhalli and Pevzner [9] have shown for the first time that the problem of sorting a *signed* permutation by reversals is polynomial: They have proved a duality theorem which equates the reversal distance with the sum of three combinatorial parameters (see Theorem 2.1 below). Based on this theorem, Hannenhalli and Pevzner proved that sorting signed permutations by reversals can be done in  $O(n^4)$  time. More recently, Berman and Hannenhalli [4] described a faster implementation that finds a minimum sequence of reversals in  $O(n^2\alpha(n))$  time, where  $\alpha()$  is the inverse

\*Department of Computer Science, Princeton University, Princeton, NJ 08544 USA and NEC Institute, Princeton, NJ. Research at Princeton University supported by the Office of Naval Research, Contract No. N00014-91-J-1463, and the NSF, Grants No. CCR-8920505 and CCR-9626862. hkl@cs.princeton.edu.

†Department of Computer Science, Sackler Faculty of Exact Sciences, Tel Aviv University, Tel-Aviv 69978 ISRAEL. Research supported in part by grants from the Ministry of Science and the Arts, Israel, and from the Israeli Academy of Sciences. shamir@math.tau.ac.il

‡Department of Computer Science, Princeton University, Princeton, NJ 08544 USA and NEC Institute, Princeton, NJ. Research at Princeton University partially supported by the NSF, Grants No. CCR-8920505 and CCR-9626862, and the Office of Naval Research, Contract No. N00014-91-J-1463. ret@cs.princeton.edu.

of Ackerman's function [1] (see also [21]).

In this study we give an  $O(n^2)$  algorithm for sorting a signed permutation of  $n$  elements, thereby improving upon the previous best known bound in [4]. In fact, if the reversal distance is  $r$ , our algorithm requires  $O(n\alpha(n)+n \times r)$  time. In the worst case  $r = \Theta(n)$ , but in the biological context it is expected that  $r \ll n$ , so the performance may be substantially better. In addition to the better time bound our study simplifies considerably both the algorithm and combinatorial structure needed for the analysis:

- The basic object we work with is an implicit representation of the overlap graph, to be defined later, in contrast with the interleaving graph in [9] and [4]. The overlap graph is combinatorially simpler than the interleaving graph. As a result, it is easier to produce a representation for the overlap graph from the input, and to maintain it while searching for reversals.
- As a consequence of our ability to work with the overlap graph we need not perform any "padding transformations", nor do we have to work with "simple permutations" as [9] and [4].
- We deal with the unoriented and oriented parts of the permutation separately, which makes the algorithm much simpler.
- The notion of a hurdle, one of the combinatorial entities defined by [9] for the duality theorem, is simplified and is addressed in a more symmetric manner.
- The search for the next reversal is much simpler, and requires no special data structures. Our algorithm computes connected components only once, and any simple implementation of it suffices to obtain the quadratic time bound. In contrast, in [4] a logarithmic number of connected components computations may be performed per reversal, using the union-find data structure.

The paper is organized as follows: Section 2 gives the necessary preliminaries. Section 3 describes how to implicitly generate and represent the overlap graph of a permutation. Sections 4 and 5 give the details of our algorithm. We summarize our results and suggest some further research in Section 6.

## 2 Preliminaries

This section gives the basic background, primarily the theory of Hannenhalli and Pevzner, on which we base our algorithm. We start with some definitions for unsigned permutations: Let  $\pi = (\pi_1, \dots, \pi_n)$  denote a permutation of  $\{1, \dots, n\}$ . Augment  $\pi$  to a permutation on  $n + 2$  vertices by adding  $\pi_0 = 0$  and  $\pi_{n+1} = n + 1$  to it. A pair  $(\pi_i, \pi_{i+1})$ ,  $0 \leq i \leq n$  is called a *gap*. Gaps are classified into two types. A gap  $(\pi_i, \pi_{i+1})$  is a

*breakpoint* of  $\pi$  if and only if  $|\pi_i - \pi_{i+1}| > 1$ , otherwise it is an *adjacency* of  $\pi$ . We denote by  $b(\pi)$  the number of breakpoints in  $\pi$ .

A *reversal*,  $\rho(i, j)$ , on a permutation  $\pi$  transforms  $\pi$  to  $\pi' = \pi\rho(i, j) = (\pi_1, \dots, \pi_{i-1}, \pi_j, \pi_{j-1}, \dots, \pi_i, \pi_{j+1}, \dots, \pi_n)$ . We say that a reversal  $\rho(i, j)$  is *acting on* the gaps  $(\pi_{i-1}, \pi_i)$  and  $(\pi_j, \pi_{j+1})$ . The *breakpoint graph*  $B(\pi)$  of a permutation  $\pi = (\pi_1, \dots, \pi_n)$  is an edge-colored graph on  $n + 2$  vertices  $\{\pi_0, \pi_1, \dots, \pi_{n+1}\} = \{0, 1, \dots, n + 1\}$ . We join vertices  $\pi_i$  and  $\pi_j$  by a *black edge* if  $(\pi_i, \pi_j)$  is a breakpoint in  $\pi$  and by a *gray edge* if  $(i, j)$  is a breakpoint in  $\pi^{-1}$ .

We define a one to one mapping  $u$  from the set of signed permutations of order  $n$  into the set of unsigned permutations of order  $2n$  as follows. Let  $\pi$  be a signed permutation. To obtain  $u(\pi)$  replace each positive element  $x$  in  $\pi$  by  $2x - 1, 2x$  and negative element  $-x$  by  $2x, 2x - 1$ . For any signed permutation  $\pi$ , let  $B(\pi) = B(u(\pi))$ . For the rest of this paper we limit the discussion to signed permutations. Note that in  $B(\pi)$  every vertex is either isolated or incident with exactly one black edge and one gray edge. Therefore, there is a unique decomposition of  $B(\pi)$  into cycles. The edges of each cycle are alternating gray and black. Call a reversal  $\rho(i, j)$  such that  $i$  is odd and  $j$  even an *even reversal*. The reversal  $\rho(2i + 1, 2j)$  on  $u(\pi)$  mimics the reversal  $\rho(i + 1, j)$  on  $\pi$ . Thus, sorting  $\pi$  by reversals is equivalent to sorting the unsigned permutation  $u(\pi)$  by even reversals. Henceforth we will consider the latter problem and by reversal we will always mean an even reversal. Let  $b(\pi) = b(u(\pi))$  and  $c(\pi)$  be the number of cycles in  $B(\pi)$ .

For an arbitrary reversal  $\rho$  on a permutation  $\pi$ , denote by  $\Delta b(\pi, \rho) = b(\pi\rho) - b(\pi)$  and  $\Delta c(\pi, \rho) = c(\pi\rho) - c(\pi)$ . When the reversal  $\rho$  and the permutation  $\pi$  will be clear from the context we will abbreviate  $\Delta b(\pi, \rho)$  to  $\Delta b$  and  $\Delta c(\pi, \rho)$  to  $\Delta c$ . The following values are taken by  $\Delta b$  and  $\Delta c$  depending upon the types of the gaps  $\rho(i, j)$  is acting on. They were first observed by Bafna and Pevzner [2], and are straightforward to verify:

1. Two adjacencies;  $\Delta c = 1$  and  $\Delta b = 2$ .
2. A breakpoint and an adjacency;  $\Delta c = 0$  and  $\Delta b = 1$ .
3. Two breakpoints each belonging to a different cycle;  $\Delta b = 0$ ,  $\Delta c = -1$ .
4. Two breakpoints of the same cycle  $C$ :
  - a.  $(\pi_i, \pi_{j+1})$  and  $(\pi_{i-1}, \pi_j)$  are gray edges;  $\Delta c = -1$ ,  $\Delta b = -2$ .
  - b. Exactly one of  $(\pi_i, \pi_{j+1})$  and  $(\pi_{i-1}, \pi_j)$  is a gray edge;  $\Delta c = 0$ ,  $\Delta b = -1$ .
  - c. Neither  $(\pi_i, \pi_{j+1})$  nor  $(\pi_{i-1}, \pi_j)$  is a gray edge, and when breaking  $C$  at  $i$  and  $j$  vertices  $i - 1$  and

$j + 1$  end up in the same path;  $\Delta b = 0, \Delta c = 0$   
 d. Neither  $(\pi_i, \pi_{j+1})$  nor  $(\pi_{i-1}, \pi_j)$  is a gray edge, and when breaking  $C$  at  $i$  and  $j$  vertices  $i - 1$  and  $j + 1$  end up in different paths;  $\Delta b = 0, \Delta c = 1$

Call a reversal *proper* if  $\Delta b - \Delta c = -1$ , i.e. it is either of type 4a, 4b, or 4d. We say that a reversal  $\rho$  is *acting on a gray edge  $e$*  if it is acting on the breakpoints which correspond to the black edges incident with  $e$ . A gray edge is *oriented* if a reversal acting on it is proper, otherwise it is *unoriented*. A cycle is *oriented* if it contains an oriented gray edge, and it is *unoriented* otherwise.

Two intervals on the real line *overlap* if their intersection is nonempty but neither properly contains the other. A graph  $G$  is an *interval overlap graph* if one can assign an interval to each vertex such that two vertices are adjacent if and only if the corresponding intervals overlap (see, e.g., [7]). For a permutation  $\pi$ , we associate with a gray edge  $(\pi_i, \pi_j)$  the interval  $[i, j]$ . The *overlap graph* of a permutation  $\pi$ , denoted  $OV(\pi)$ , is the interval overlap graph of the gray edges of  $B(\pi)$ . Namely, the vertex set of  $OV(\pi)$  is the set of gray edges in  $B(\pi)$ , and two vertices are connected if the intervals associated with their gray edges overlap. Throughout this paper whenever we talk about the *representation of  $OV(\pi)$*  we refer to this canonical representation. Note that all the endpoints of intervals in this representation are distinct integers. We also identify a vertex in  $OV(\pi)$  with the edge it represents and with its interval in the representation. Thus, the endpoints of a gray edge are actually the endpoints of the interval representing the corresponding vertex in  $OV(\pi)$ . A connected component of  $OV(\pi)$  that contains an oriented edge is called an *oriented component*, otherwise it is called an *unoriented component*.

Let  $X$  be a set of gray edges in  $B(\pi)$ . Define  $\min(X) = \min\{i \mid (\pi_i, \pi_j) \in X\}$ ,  $\max(X) = \max\{j \mid (\pi_i, \pi_j) \in X\}$  and  $\text{span}(X) = [\min(X), \max(X)]$ . Equivalently, one can look at the interval overlap representation of  $OV(\pi)$  mentioned above and define the span of a set of vertices  $X$  as the minimum interval which contains all the intervals of vertices in  $X$ .

The major object our algorithm will work with is  $OV(\pi)$  though for efficiency considerations we will avoid generating it explicitly. In contrast, Pevzner and Hannenhalli worked with the *interleaving graph  $H_\pi$* , whose vertices are the alternating cycles of  $B(\pi)$  and two cycles  $C_1$  and  $C_2$  are connected by an edge in  $H_\pi$  iff there exists a gray edge  $e_1 \in C_1$  and a gray edge  $e_2 \in C_2$  that overlap.

The following lemma and its corollary imply that the partition imposed by the connected components of

$OV(\pi)$  on the set of gray edges is identical to the one imposed by the connected components of  $H_\pi$ :

LEMMA 2.1. *Assume  $M$  is a set of gray edges in  $B(\pi)$  that corresponds to a connected component in  $OV(\pi)$  then  $\min(M)$  is even and  $\max(M)$  is odd.*

*Proof.* Assume  $\min(M)$  is odd, then  $\pi_{\min(M)} + 1$  and  $\pi_{\min(M)} - 1$  must both be in  $\text{span}(M)$  (i.e. there exist  $l_1, l_2 \in \text{span}(M)$  such that  $\pi_{l_1} = \pi_{\min(M)} + 1$  and  $\pi_{l_2} = \pi_{\min(M)} - 1$ ). Thus,  $\pi_{\min(M)}$  is neither the maximum nor the minimum element in the set  $\{\pi_i \mid i \in \text{span}(M)\}$ . Hence, either the maximum element or the minimum element in  $\text{span}(M)$  is  $\pi_j$  for some  $\min(M) < j < \max(M)$ . By the definition of  $B(\pi)$  there must be a gray edge  $(\pi_j, \pi_l)$  for some  $l \notin \text{span}(M)$ , contradicting the fact that  $M$  is a connected component in  $OV(\pi)$ . The proof that  $\max(M)$  is odd is similar. ■

COROLLARY 2.1. *Every connected component of  $OV(\pi)$  corresponds to the set of gray edges of a union of cycles.*

*Proof.* Assume, by contradiction, that  $C$  is a cycle whose gray edges belong to at least two connected components in  $OV(\pi)$ . Assume  $M_1$  and  $M_2$  are two of these components such that there are two consecutive gray edges  $e_1 \in M_1$  and  $e_2 \in M_2$  along  $C$ . Since the spans of different connected components in  $OV(\pi)$  cannot overlap there are two different cases to consider. 1.  $\text{span}(M_2) \subseteq \text{span}(M_1)$  (the case  $\text{span}(M_1) \subseteq \text{span}(M_2)$  is symmetric). Since  $e_1$  and  $e_2$  are in different components they cannot overlap. Thus, either the right endpoint of  $e_2$  is even and equals  $\max(M_2)$  or the left endpoint of  $e_2$  is odd and equals  $\min(M_2)$ . In both cases we obtained a contradiction to Lemma 2.1. 2.  $\text{span}(M_2)$  and  $\text{span}(M_1)$  are disjoint intervals. W.l.o.g. assume that  $\max(M_1) < \min(M_2)$ . The right endpoint of  $e_1$  is even and equals  $\max(M_1)$ , and that contradicts Lemma 2.1. ■

Let  $\pi_{i_1}, \pi_{i_2}, \dots, \pi_{i_k}$  be the subsequence of  $0, \pi_1, \dots, \pi_n, n + 1$  consisting of those elements incident with gray edges in  $B(\pi)$  that occur in unoriented components of  $OV(\pi)$ . Order  $\pi_{i_1}, \pi_{i_2}, \dots, \pi_{i_k}$  on a circle  $CR$  such that  $\pi_{i_j}$  follows  $\pi_{i_{j-1}}$  for  $2 \leq j \leq k$  and  $\pi_{i_1}$  follows  $\pi_{i_k}$ . Let  $M$  be an unoriented connected component in  $OV(\pi)$ . Let  $E(M) \subset \{\pi_{i_1}, \pi_{i_2}, \dots, \pi_{i_k}\}$  be the set of endpoints of the edges in  $M$ . An unoriented component  $M$  is a *hurdle* if the elements of  $E(M)$  occur consecutively on  $CR$ .

This definition of a hurdle is different from the one given by Hannenhalli and Pevzner [9]. It is simpler in the sense that minimal hurdles and the maximal one do not have to be treated in a different way. Using Corollary 2.1 above one can prove that the hurdles as we defined them are identical to the ones defined by

Hannenhalli and Pevzner. Let  $h(\pi)$  denote the number of hurdles in a permutation  $\pi$ .

A hurdle is *simple* if when one deletes it from  $OV(\pi)$  no other unoriented component becomes a hurdle, and it is a *super hurdle* otherwise. A *fortress* is a permutation with an odd number of hurdles all of which are super hurdles.

The following theorem was proved by Hannenhalli and Pevzner.

**THEOREM 2.1.** [9] *The minimum number of reversals required to sort a permutation  $\pi$  is  $b(\pi) - c(\pi) + h(\pi)$  unless  $\pi$  is a fortress in which case exactly one additional reversal is necessary and sufficient.*

Denote by  $d(\pi)$  the reversal distance of  $\pi$ . I.e.  $d(\pi) = b(\pi) - c(\pi) + h(\pi) + 1$  if  $\pi$  is a fortress and  $d(\pi) = b(\pi) - c(\pi) + h(\pi)$  otherwise. Following the theory developed in [9] it turns out that given a permutation  $\pi$  with  $h(\pi) > 0$  one can perform  $t = \lceil h(\pi)/2 \rceil$  reversals and transform  $\pi$  to a permutation  $\pi'$  such that  $h(\pi') = 0$  and  $d(\pi') = d(\pi) - t$ . Our method of “clearing the hurdles” uses the same theory developed by Hannenhalli and Pevzner. In Section 5 we describe an efficient implementation of this process which uses the implicit representation of the overlap graph  $OV(\pi)$ . Our implementation runs in  $O(n)$ -time assuming  $OV(\pi)$  is already partitioned into its connected components. Recently, Berman and Hannenhalli [4] gave an  $O(n\alpha(n))$  algorithm for computing the connected components of an interval overlap graph given implicitly by its representation. Using their algorithm we can clear the hurdles from a permutation in  $O(n\alpha(n))$  time.

The overlap graph of  $\pi'$ ,  $OV(\pi')$ , has only oriented components. In Section 4 we prove that in the neighborhood of any oriented gray edge  $e$  there is an oriented gray edge  $e_1$  ( $e_1$  could be the same as  $e$ ) such that a reversal acting on  $e_1$  does not create new hurdles. Call such a reversal a *safe reversal*. We develop an efficient algorithm to locate a safe reversal in a permutation with at least one oriented gray edge. Our algorithm uses only an implicit representation of the overlap graph and it runs in  $O(n)$  time.

### 3 Representing the overlap graph

We assume that the input is given as a sequence of  $n$  signed integers representing  $\pi^0$ . First the permutation  $\pi = u(\pi^0)$  is constructed as described in the previous section. The second stage is to construct an interval overlap representation of  $OV(\pi)$ . By that we mean a listing of the intervals corresponding to the gray edges which are the vertices in  $OV(\pi)$  ordered in increasing left endpoint order.

One could construct the representation in linear time as follows. Construct an array  $b$  representing  $\pi^{-1}$ .

For each  $i$  from 0 to  $2n$  do: if  $|b(i) - b(i+1)| > 1$  then add the interval  $[b(i), b(i+1)]$  (or  $[b(i+1), b(i)]$ , if  $b(i) > b(i+1)$ ) to the representation.

### 4 Eliminating oriented components

First we introduce some notation. Recall that the vertices of  $OV(\pi)$  are the gray edges of  $B(\pi)$ . In order to avoid confusion we will usually refer to them as vertices of  $OV(\pi)$ . Hence a vertex of  $OV(\pi)$  is *oriented* if the corresponding gray edge is oriented and it is *unoriented* otherwise. Let  $e$  be a vertex in  $OV(\pi)$ . Denote by  $r(e)$  the reversal acting on the gray edge corresponding to  $e$ . Denote by  $N(e)$  the set of neighbors of  $e$  in  $OV(\pi)$  including  $e$  itself. Denote by  $ON(e)$  the subset of  $N(e)$  containing the oriented vertices and by  $UN(e)$  the subset of  $N(e)$  containing the unoriented vertices.

In this section we prove that if an oriented vertex  $e$  exists in  $OV(\pi)$  then there exists an oriented vertex  $f \in ON(e)$  such that  $r(f)$  is proper and safe. We also describe an algorithm that finds a proper safe reversal in a permutation that contains at least one oriented edge.

We start with the following useful observation:

**OBSERVATION 4.1.** *Let  $e$  be a vertex in  $OV(\pi)$  and let  $\pi' = \pi r(e)$ .  $OV(\pi')$  could be obtained from  $OV(\pi)$  by the following operations. 1) Complement the graph induced by  $OV(\pi)$  on  $N(e)$ , and flip the orientation of every vertex in  $N(e)$ . 2) If  $e$  is oriented in  $OV(\pi)$  then remove it from  $OV(\pi)$ . 3) If there exists an oriented edge  $e'$  in  $OV(\pi)$  with  $r(e) = r(e')$  then remove  $e'$  from  $OV(\pi)$ .*

Note that if  $e$  is an oriented vertex in a component  $M$  of  $OV(\pi)$ ,  $M - \{e\}$  may split into several components in  $OV(\pi')$ . Denote these components  $M'_1(e), \dots, M'_k(e)$ , where  $k \geq 1$ . We will refer to  $M'_i(e)$  simply as  $M'_i$  whenever  $e$  is clear from the context.

Let  $C$  be a clique of oriented vertices in  $OV(\pi)$ . We say that  $C$  is *happy* if for every oriented vertex  $e \notin C$  and every vertex  $f \in C$  such that  $(e, f) \in E(OV(\pi))$  there exists an oriented vertex  $g \notin C$  such that  $(g, e) \in E(OV(\pi))$  and  $(g, f) \notin E(OV(\pi))$ . Our first theorem claims that one of vertices in any happy clique defines a safe proper reversal.

**THEOREM 4.1.** *Let  $C$  be a happy clique and let  $e$  be a vertex in  $C$  such that  $|UN(e')| \leq |UN(e)|$  for every  $e' \in C$  then the reversal  $r(e)$  is safe.*

*Proof.* Let  $\pi' = \pi r(e)$  and assume by contradiction that  $M'_i(e)$  is unoriented for some  $1 \leq i \leq k$ . Clearly  $N(e) \cap M'_i \neq \emptyset$ .

Assume there exists  $y \in N(e) \cap M'_i$  such that  $y \notin C$ . Clearly  $y$  must be oriented in  $OV(\pi)$  and since  $C$  is happy it must also have an oriented neighbor  $y'$  such that  $(y', e) \notin E(OV(\pi))$ . Since  $y'$  is not adjacent to  $e$  in  $OV(\pi)$  it stays oriented and adjacent to  $y$  in

$OV(\pi')$  in contradiction with the assumption that  $M'_i$  is unoriented. Hence we may assume that  $N(e) \cap M'_i \subseteq C$ .

Let  $y \in N(e) \cap M'_i$  and let  $z \in UN(e)$ . Vertex  $z$  is oriented in  $OV(\pi')$  and if it is connected to  $y$  in  $OV(\pi')$  we obtain a contradiction. Hence,  $z$  and  $y$  are not connected in  $OV(\pi')$ , so they must be connected in  $OV(\pi)$ . Hence we obtain that  $UN(e) \subseteq UN(y)$  in  $OV(\pi)$ . Corollary 2.1 implies that component  $M'_i$  cannot contain  $y$  alone. Thus  $y$  must have a neighbor  $x$  in  $M'_i$ . Since  $N(e) \cap M'_i \subseteq C$  vertex  $x$  is not adjacent to  $e$  in  $OV(\pi)$ . So we obtain that  $(x, y) \in OV(\pi)$ ,  $(x, e) \notin OV(\pi)$ , and  $x$  is unoriented in  $OV(\pi)$ . Since we already proved that  $UN(e) \subseteq UN(y)$  this implies that  $UN(e) \subset UN(y)$  in contradiction with the choice of  $e$ . ■

**THEOREM 4.2.** *Let  $e$  be an oriented vertex in a component  $M$  of  $OV(\pi)$ . There exists an oriented vertex  $f \in ON(e)$  such that  $M'_1(f), \dots, M'_k(f)$  are all oriented in  $OV(\pi')$ , where  $\pi' = \pi r(f)$ .*

*Proof.* By Theorem 4.1 it suffices to show that there exists a happy clique  $C$  in  $ON(e)$ .

Let  $Ext(e) = \{x \in ON(e) \mid \text{there exists } y \in ON(x) \text{ such that } y \notin ON(e)\}$ . I.e.  $Ext(e)$  contains all oriented neighbors of  $e$  which have oriented neighbors outside of  $ON(e)$ .

Case 1:  $Ext(e) = ON(e) - \{e\}$ . Set  $C = \{e\}$ .

Case 2:  $Ext(e) \subset ON(e) - \{e\}$ . let  $D^0 = ON(e) - Ext(e)$ . For  $j \geq 0$ , while  $D^j$  is not a clique let  $K^j$  be a maximal clique in  $D^j$  and define  $D^{j+1} = D^j - K^j$ . Let  $D^k$ ,  $k \geq 0$  be the final clique and set  $C = D^k$ .

It is straightforward to verify that in each of the two cases  $C$  is indeed a happy clique. ■

Though  $OV(\pi)$  has at most  $n + 1$  vertices, it may have a superlinear number of edges. Thus creating and maintaining it while looking for reversals may be expensive. In contrast, the interval overlap representation of  $OV(\pi)$  has linear size. In the next section we describe an algorithm that will find an oriented edge  $e$  such that  $r(e)$  is safe given a representation of  $OV(\pi)$ . The algorithm first finds a happy clique  $C$  and then searches for the vertex with maximum unoriented degree in  $C$ . According to Theorem 4.1 that vertex defines a safe reversal. The time complexity of each stage is  $O(n)$ .

**4.1 Finding a happy clique** In this section we give an algorithm that locates a happy clique in  $OV(\pi)$ . Let  $e_1, \dots, e_k$  be the oriented vertices in  $OV(\pi)$  in increasing left endpoint order. Let  $L(e)$  and  $R(e)$  be the left and right endpoints, respectively, of vertex  $e$  in the realization of  $OV(\pi)$ . The algorithm traverses the oriented vertices in  $OV(\pi)$  according to this order. After traversing  $e_1, \dots, e_i$ ,  $1 \leq i \leq k$  the algorithm maintains a happy clique  $C_i$  in the subgraph of  $OV(\pi)$

induced by these vertices. Assume  $|C_i| = j$ ,  $j \leq i$  and let  $e_{i_1}, \dots, e_{i_j}$  be the vertices in  $C_i$  where  $i_1 < i_2 < \dots < i_j$ . The vertices of  $C_i$  are maintained in a linked list ordered in increasing left endpoint order. If there exists an interval that contains all the intervals in  $C_i$  then the algorithm maintains one such interval  $t_i$ . The clique  $C_i$  and the vertex  $t_i$  (if exists) satisfy the following invariant.

**INVARIANT 4.1.**

- 1) Every vertex  $e_l \notin C_i$ ,  $l \leq i$ , such that  $L(e_{i_1}) < L(e_l)$  must be adjacent to  $t_i$ , i.e.,  $R(e_l) > R(t_i)$ .
- 2) Every vertex  $e_l \notin C_i$ ,  $L(e_l) < L(e_{i_1})$  that is adjacent to a vertex in  $C_i$  is either adjacent to an interval  $e_p$  such that  $R(e_p) < L(e_{i_1})$  or adjacent to  $t_i$ .

The fact that  $C_i$  is happy in the subgraph induced by  $e_1, \dots, e_i$  follows from this invariant. We initialize the algorithm by setting  $C_1 = \{e_1\}$ . Initially,  $t_1$  is not defined. Let the current interval be  $e_{i+1}$ . If  $R(e_{i_j}) < L(e_{i+1})$  then  $C_i$  is guaranteed to be happy in  $OV(\pi)$  since all remaining oriented vertices are not adjacent to  $C_i$ . Hence the algorithm stops and returns  $C_i$  as the answer.

We now assume that  $R(e_{i_j}) \geq L(e_{i+1})$  and show how to obtain  $C_{i+1}$  and  $t_{i+1}$ . We have to consider the following cases.

Case 1. The interval  $t_i$  is defined and  $R(e_{i+1}) > R(t_i)$ . Continue with  $C_{i+1} := C_i$  and  $t_{i+1} := t_i$ .

Case 2. The interval  $t_i$  is not defined or  $R(e_{i+1}) \leq R(t_i)$ .

- a)  $R(e_{i+1}) > R(e_{i_j})$  and  $L(e_{i+1}) \leq R(e_{i_1})$ .  $C_{i+1}$  is obtained by adding  $e_{i+1}$  to  $C_i$  and  $t_{i+1} := t_i$ .

- b)  $R(e_{i+1}) > R(e_{i_j})$  and  $L(e_{i+1}) > R(e_{i_1})$ . The clique  $\overline{C}_{i+1}$  consists of  $e_{i+1}$  alone and  $t_{i+1} := t_i$ .

- c)  $R(e_{i+1}) < R(e_{i_j})$ . As in the previous case  $C_{i+1} = \overline{\{e_{i+1}\}}$ . In this case  $t_{i+1}$  is set to  $e_{i_j}$ , the last interval in  $C_i$ .

The following theorem proves that the algorithm we described produces a happy clique.

**THEOREM 4.3.** *Let  $C_l$  be the current clique when the algorithm stops. Then  $C_l$  is a happy clique in  $OV(\pi)$ .*

*Proof.* A straightforward induction on the number of oriented vertices traversed by the algorithm proves that  $C_l$  and  $t_l$  satisfy Invariant 4.1.

The algorithm stops either when  $R(e_{i_j}) < L(e_{i+1})$  or when  $l = k$  where  $k$  is the number of oriented vertices. In either case since  $C_l$  is happy in the subgraph induced by  $e_1, \dots, e_l$  it must be happy in  $OV(\pi)$ . ■

The running time of the algorithm is proportional to the number of oriented vertices traversed since a constant amount of work is performed per such vertex.

**4.2 Searching the happy clique**

After locating a happy clique  $C$  in  $OV(\pi)$  we need to search it for a vertex with a maximum number of unoriented neighbors. In this section we give an algorithm that performs this task.

Let  $e_1, \dots, e_j$  be the intervals in  $C$  ordered in increasing left endpoint order. Clearly,  $L(1) < L(2) < \dots < L(j) < R(1) < R(2) < \dots < R(j)$ . Thus the endpoints of the  $j$  vertices in  $C$  partition the line into  $2j + 1$  disjoint intervals  $I_0, \dots, I_{2j}$ , where  $I_0 = (-\infty, L(1))$ ,  $I_l = (L(l), L(l + 1)]$  for  $1 \leq l < j$ ,  $I_j = (L(j), R(1))$ ,  $I_l = (R(l - j), R(l - j + 1)]$  for  $j < l < 2j$  and  $I_{2j} = (R(j), \infty)$ . The algorithm consists of the following three stages.

**Stage 1:** Let  $e$  be an unoriented vertex that overlaps with the interval  $[L(1), R(j)]$ . Mark each of  $e$ 's endpoints with the index of the interval that contains it.

**Stage 2:** Let  $o$  be an array of  $j$  counters each corresponding to a vertex in  $C$ . The intention is to assign values to  $o$  such that the sum  $\sum_{i=1}^j o[i]$  is the unoriented degree of the vertex  $e_l \in C$ . The counters are initialized to zero. For each unoriented vertex  $e$  that overlaps with the interval  $[L(1), R(j)]$  we change at most four of the counters as follows. Let  $I_l$  and  $I_r$  be the intervals in which  $L(e)$  and  $R(e)$  occur respectively. We may assume  $l < r$  as otherwise  $e$  is not adjacent to any vertex in  $C$  and we can ignore it. We continue according to one of the following cases.

**Case 1:**  $r \leq j$ . All the vertices from  $e_{l+1}$  to  $e_r$  are adjacent to  $e$  hence we increment  $o[l + 1]$  and decrement  $o[r + 1]$  (if  $r < j$ ).

**Case 2:**  $j \leq l$ . All the vertices from  $e_{l-j+1}$  to  $e_{r-j}$  are adjacent to  $e$  hence we increment  $o[l - j + 1]$  and decrement  $o[r - j]$ .

**Case 3:**  $l < j$  and  $j < r$ . Let  $m = \min\{l, r - j\}$ . If  $m > 0$  then all the vertices from  $e_1$  to  $e_m$  are adjacent to  $e$  hence we increment  $o[1]$  and decrement  $o[m + 1]$ . Similarly let  $M = \max\{l, r - j\}$ . If  $M < j$  then the vertices from  $e_{l+1}$  to  $e_j$  intersect  $e$  hence we increment the counter  $o[l + 1]$ .

**Stage 3:** Compute  $f = \max_l \{\sum_{i=1}^j o[i] \mid 1 \leq l \leq j\}$ . Return  $e_f$ .

The following theorem summarizes the result of this section. We omit the proof which is straightforward.

**THEOREM 4.4.** *Given a clique  $C$ , the vertex  $e_f \in C$  computed by the algorithm above has maximum unoriented degree among the vertices in  $C$ .*

The complexity of the algorithm is proportional to the size of  $C$  plus the number of unoriented vertices in  $OV(\pi)$ . That is  $O(n)$ .

## 5 Clearing the hurdles

In case there are unoriented components in  $OV(\pi)$  there exists a sequence  $r_1, \dots, r_t$  of  $t$  reversals that transform  $\pi$  into  $\pi'$  such that  $d(\pi') = d(\pi) - t$ , where  $t = \lceil h(\pi)/2 \rceil$ . In this section we summarize the characterization given by Hannenhalli and Pevzner for these  $t$  reversals and outline how to find them using our implicit representation of  $OV(\pi)$ .

We will use the following definitions. A reversal merges hurdles  $H_1$  and  $H_2$  if it acts on two breakpoints one incident with a gray edge in  $H_1$  and the other incident with a gray edge in  $H_2$ . Recall the circle  $CR$  defined in Section 2 of the endpoints of the edges in the unoriented components of  $OV(\pi)$  is ordered consistently with their order in  $\pi$ . Two hurdles  $H_1$  and  $H_2$  are consecutive if their sets of endpoints  $E(H_1)$  and  $E(H_2)$  occur consecutively on  $CR$ . I.e. there is no hurdle  $H$  such that  $E(H)$  separates  $E(H_1)$  and  $E(H_2)$  on  $CR$ .

The following lemmas were essentially proved by Hannenhalli and Pevzner though stated differently in their paper.

**LEMMA 5.1.** ([9]) *Let  $\pi$  be a permutation with an even number, say  $2k$ , of hurdles. Any sequence of  $k - 1$  reversals each of which merges two non-consecutive hurdles followed by a reversal merging the remaining two hurdles will transform  $\pi$  into  $\pi'$  such that  $d(\pi') = d(\pi) - k$  and  $\pi'$  has only oriented components.*

**LEMMA 5.2.** ([9]) *Let  $\pi$  be a permutation with an odd number, say  $2k + 1$ , of hurdles. If at least one hurdle  $H$  is simple then a reversal acting on two breakpoints incident with edges in  $H$  transforms  $\pi$  into  $\pi'$  with  $2k$  hurdles such that  $d(\pi') = d(\pi) - 1$ . If  $\pi$  is a fortress then a sequence of  $k - 1$  reversals merging pairs of non-consecutive hurdles followed by two additional merges of pairs of consecutive hurdles (one merges two original hurdles and the next merges a hurdle created by the first and the last original hurdle) will transform  $\pi$  into  $\pi'$  such that  $d(\pi') = d(\pi) - (k + 1)$  and  $\pi'$  has only oriented components.*

We now outline how to turn these lemmas into an algorithm that finds a particular sequence of reversals  $r_1, \dots, r_t$  with the properties described above. First  $OV(\pi)$  is decomposed into connected components as described in [4]. One then has to identify those unoriented components that are hurdles. This task could be done by traversing the endpoints of the circle  $CR$  counting the number of elements in each run of consecutive endpoints belonging to the same component. If a run contains all endpoints of a particular unoriented component  $M$  then  $M$  is an hurdle.

In a similar fashion one could check for each hurdle whether it is a simple hurdle or a super hurdle. While

traversing the cycle a list of the hurdles in the order they occur on  $CR$  should be created. At the next stage this list would be used to identify correct hurdles to merge.

We assume that given an endpoint one can locate its connected component in constant time. It is easy to verify that the data could be maintained such that it is possible to do so.

**THEOREM 5.1.** *Given  $OV(\pi)$  decomposed into its connected components, the algorithm outlined above finds  $t$  reversals such that when we apply them to  $\pi$  we obtain  $\pi'$  which is hurdle-free and  $d(\pi') = d(\pi) - t$ . It could be implemented to run in  $O(n)$ -time.*

*Proof.* The correctness follows from Lemma 5.1 and 5.2. The time bound is achieved if we always merge hurdles that are separated by a single hurdle. If the  $i$ th merge merged hurdles  $H_1$  and  $H_2$  that are separated by  $H$  then  $H$  should be merged in the  $i + 1$ st merge. Carrying out the merges that way guarantees that the span of each hurdle  $H$  overlaps at most two merging reversals the second of which eliminates  $H$ . ■

## 6 Summary

Figure 1 gives a schematic description of the algorithm.

```

algorithm SIGNED REVERSALS( $\pi$ );
/*  $\pi$  is a signed permutation */
1. Compute the connected components of  $OV(\pi)$ .
2. Clear the hurdles.
3. while  $\pi$  is not sorted do :
/* iteration */
begin
  a. find a happy clique  $C$  in  $OV(\pi)$ .
  b. find a vertex  $e_f \in C$  with maximum unoriented
     degree, and perform a safe reversal on  $e_f$ ;
  c. update  $\pi$  and the representation of  $OV(\pi)$ .
end
4. output the sequence of reversals.

```

Figure 1: An algorithm for sorting signed permutations

### THEOREM 6.1.

*Algorithm SIGNED REVERSALS finds the reversal distance  $r$  in  $O(n\alpha(n) + r \times n)$  time, and in particular in  $O(n^2)$  time.*

*Proof.* The correctness of the algorithm follows from Theorem 2.1, Theorem 4.1 and Lemmas 5.1 and 5.2.

Step 1 takes  $O(n\alpha(n))$  by the algorithm of Berman and Hannenhalli [4]. Step 2 takes  $O(n)$  time by Theorem 5.1. Step 3 takes  $O(n)$  time per reversal, by the discussion in Section 4. ■

It is an intriguing open question whether a faster algorithm for sorting signed permutations by reversals exists. It certainly might be the case that one can find an optimal sequence of reversals much faster. To date, no nontrivial lower bound is known for this problem.

## Acknowledgments

We thank Sridhar Hannenhalli, Pavel Pevzner and Izik Pe'er for their comments on a preliminary version of this paper.

## References

- [1] W. Ackermann. Zum hilbertschen aufbau der reellen zahlen. *Math. Ann.*, 99:118–133, 1928.
- [2] V. Bafna and P. Pevzner. Genome rearrangements and sorting by reversals. In *Proc. 34th IEEE Symp. of the Foundations of Computer Science*, pages 148–157. IEEE Computer Society Press, 1994. To appear in SIAM J. of Computing.
- [3] V. Bafna and P. Pevzner. Sorting permutations by transpositions. In *Proceedings of the 6th Annual Symposium on Discrete Algorithms*, pages 614–623. ACM Press, Jan. 1995.
- [4] P. Berman and S. Hannenhalli. Fast sorting by reversal. In *Proc. Combinatorial Pattern Matching (CPM) 1996*, 1996.
- [5] A. Caprara. Sorting by reversals is difficult. Technical report, DEIS, University of Bologna, April 1996.
- [6] T. Dobzhansky and A. H. Sturtevant. Inversions in the chromosomes of *drosophila pseudoobscura*. *Genetics*, 23:28–64, 1938.
- [7] M. C. Golumbic. *Algorithmic Graph Theory and Perfect Graphs*. Academic Press, New York, 1980.
- [8] S. Hannenhalli. Polynomial algorithm for computing translocation distance between genomes. Technical Report CSE-95-005, Pennsylvania State University, 1995.
- [9] S. Hannenhalli and P. Pevzner. Transforming cabbage into turnip (polynomial algorithm for sorting signed permutations by reversals). In *Proceedings of the Twenty-Seventh Annual ACM Symposium on Theory of Computing*, pages 178–189, Las Vegas, Nevada, 29 May–1 June 1995.
- [10] S. Hannenhalli and P. Pevzner. Transforming men into mice (polynomial algorithm for genomic distance problems). In *Proc. IEEE Symp. of the Foundations of Computer Science*, 1995.
- [11] S. B. Hoot and J. D. Palmer. Structural rearrangements, including parallel inversions, within the chloroplast genome of *Anemone* and related genera. *J. Molecular Evolution*, 38:274–281, 1994.
- [12] J. Kececioglu and R. Ravi. Physical mapping of chromosomes using unique probes. In *Proc. sixth annual ACM-SIAM Symp. on Discrete Algorithms (SODA 95)*, pages 604–613. ACM Press, 1995.
- [13] J. Kececioglu and D. Sankoff. Efficient bounds for oriented chromosome inversion distance. In *Proc. of*

*5th Ann. Symp. on Combinatorial Pattern Matching*, pages 307–325. Springer, 1994. LNCS 807.

- [14] J. Kececioglu and D. Sankoff. Exact and approximation algorithms for sorting by reversals, with application to genome rearrangement. *Algorithmica*, 13(1/2):180–210, Jan. 1995.
- [15] J. D. Palmer and L. A. Herbon. Tricircular mitochondrial genomes of Brassica and Raphanus: reversal of repeat configurations by inversion. *Nucleic Acids Research*, 14:9755–9764, 1986.
- [16] J. D. Palmer and L. A. Herbon. Unicircular structure of the Brassica hirta mitochondrial genome. *Current Genetics*, 11:565–570, 1987.
- [17] J. D. Palmer and L. A. Herbon. Plant mitochondrial DNA evolves rapidly in structure, but slowly in sequence. *J. Molecular Evolution*, 28:87–97, 1988.
- [18] J. D. Palmer, B. Osorio, and W. Thompson. Evolutionary significance of inversions in legume chloroplast DNAs. *Current Genetics*, 14:65–74, 1988.
- [19] D. Sankoff. Edit distance for genome comparison based on non-local operations. *Lecture Notes in Computer Science*, 644:121–135, 1992.
- [20] D. Sankoff, R. Cedergren, and Y. Abel. Genomic divergence through gene rearrangement. *Methods in Enzymology*, 183:428–438, 1990.
- [21] R. E. Tarjan. Efficiency of a good but not linear set union algorithm. *J. ACM*, 22(2):215–225, 1979.