

## Exact and Approximation Algorithms for Sorting by Reversals, with Application to Genome Rearrangement

J. Kececioglu<sup>1</sup> and D. Sankoff<sup>2</sup>

**Abstract.** Motivated by the problem in computational biology of reconstructing the series of chromosome inversions by which one organism evolved from another, we consider the problem of computing the shortest series of reversals that transform one permutation to another. The permutations describe the order of genes on corresponding chromosomes, and a *reversal* takes an arbitrary substring of elements, and reverses their order.

For this problem, we develop two algorithms: a greedy approximation algorithm, that finds a solution provably close to optimal in  $O(n^2)$  time and  $O(n)$  space for  $n$ -element permutations, and a branch-and-bound exact algorithm, that finds an optimal solution in  $O(mL(n, n))$  time and  $O(n^2)$  space, where  $m$  is the size of the branch-and-bound search tree, and  $L(n, n)$  is the time to solve a linear program of  $n$  variables and  $n$  constraints. The greedy algorithm is the first to come within a constant factor of the optimum; it guarantees a solution that uses no more than twice the minimum number of reversals. The lower and upper bounds of the branch-and-bound algorithm are a novel application of maximum-weight matchings, shortest paths, and linear programming.

In a series of experiments, we study the performance of an implementation on random permutations, and permutations generated by random reversals. For permutations differing by  $k$  random reversals, we find that the average upper bound on reversal distance estimates  $k$  to within one reversal for  $k < \frac{1}{2}n$  and  $n \leq 100$ . For the difficult case of random permutations, we find that the average difference between the upper and lower bounds is less than three reversals for  $n \leq 50$ . Due to the tightness of these bounds, we can solve, to optimality, problems on 30 elements in a few minutes of computer time. This approaches the scale of mitochondrial genomes.

**Key Words.** Computational biology, Approximation algorithms, Branch-and-bound algorithms, Experimental analysis of algorithms, Edit distance, Permutations, Sorting by reversals, Chromosome inversions, Genome rearrangements.

**1. Introduction.** Much research has been devoted to efficient algorithms for the edit distance between two strings, that is, the minimum number of insertions, deletions, and substitutions to transform one string into another. Motivation

---

<sup>1</sup> Department of Computer Science, The University of Georgia, Athens, GA 30602, USA. kece@cs.uga.edu. This research was supported by a postdoctoral fellowship from the Program in Mathematics and Molecular Biology of the University of California at Berkeley under National Science Foundation Grant DMS-8720208, and by a fellowship from the Centre de recherches mathématiques of the Université de Montréal.

<sup>2</sup> Centre de recherches mathématiques, Université de Montréal, Case postal 6128, succursale A, Montréal, Québec H3C 3J7, Canada. sankoff@ere.umontreal.ca. This research was supported by grants from the Natural Sciences and Engineering Research Council of Canada, and the Fonds pour la formation de chercheurs et l'aide à la recherche (Québec). The author is a fellow of the Canadian Institute for Advanced Research.

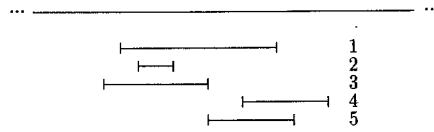


Fig. 1. Evolution of the pea chloroplast genome by five overlapping inversions.

comes in large part from computational biology: at the level of individual characters, genetic sequences mutate by these operations, so edit distance is a useful measure of evolutionary distance.

At the chromosome level, however, genetic sequences mutate by more global *genome rearrangements*, such as the reversal of a substring (*inversion*), the deletion and subsequent reinsertion of a substring far from its original site (*transposition*), the copying of a substring (*duplication*), and the exchange of prefixes or suffixes of two chromosomes in the same organism (*translocation*). An inversion, which takes a substring of unrestricted size and replaces it by its reverse in one operation, has the effect of reversing the order of the genes contained within the substring, and is perhaps the most common of these operations [21, pp. 174–175], especially in organisms with one chromosome.

For example, the only major difference between the gene orders of two of the most well-known bacteria, *Escherichia coli* and *Salmonella typhimurium*, is an inversion of a long substring of the chromosomal sequence [17]. In plants, Palmer *et al.* [18] modeled the evolution of part of the pea chloroplast genome, which is also a single chromosome, in terms of five successive overlapping inversions, as we illustrate in Figure 1. In the fruit fly, genus *Drosophila*, inversions are a far more frequent reflection of differences between and within species than translocation or other processes [4, p. 155].

The importance of inversion in these examples suggests that algorithmic study of genome rearrangement by inversion alone is a worthwhile step in the study of evolutionary distance at the level of the chromosome. Once this problem is understood, other processes such as transposition [19] and translocation [16] can be added to refine the model.

In the mathematical problem that we consider, we are given the order of  $n$  genes in two related single-chromosome organisms or two related organelles, which we represent by permutations  $\sigma = (\sigma_1 \ \sigma_2 \ \cdots \ \sigma_n)$  and  $\tau = (\tau_1 \ \tau_2 \ \cdots \ \tau_n)$ .<sup>3</sup> (In this notation,  $\sigma_i$  denotes  $\sigma(i)$ .) Such gene orders often come from genetic *maps*, that are the distillation of the work of many experimental geneticists. In current practice, the positions of the genes are increasingly found by sequence comparison, or DNA hybridization, as opposed to the mapping experiments of traditional genetics.

<sup>3</sup> Genes in one organism may be missing in the other. We assume, however, that such genes can be removed from the analysis, and that gene insertions and deletions can be analyzed separately.

We model an inversion by the reversal of an interval of elements. Formally, a *reversal* of interval  $[i, j]$  is the permutation<sup>4</sup>

$$\rho = \begin{pmatrix} i & i+1 & \cdots & j \\ j & j-1 & \cdots & i \end{pmatrix}.$$

Applying  $\rho$  to  $\sigma$  by the composition<sup>5</sup>  $\sigma \cdot \rho$  has the effect of reversing the order of genes  $\sigma_i, \sigma_{i+1}, \dots, \sigma_j$ . Our problem is the following.

**DEFINITION.** The *reversal distance problem* on permutations is, given permutations  $\sigma$  and  $\tau$ , find a series of reversals  $\rho_1, \rho_2, \dots, \rho_d$  such that

$$\sigma \cdot \rho_1 \cdot \rho_2 \cdots \rho_d = \tau,$$

and  $d$  is minimum.

We call  $d$  the *reversal distance*<sup>6</sup> between  $\sigma$  and  $\tau$ . Like edit distance, it satisfies the axioms of a metric. Reversal distance measures the amount of evolution that must have taken place at the chromosome level, assuming evolution proceeded by inversion.

Notice that the reversal distance between  $\sigma$  and  $\tau$  is equal to the reversal distance between  $\tau^{-1} \cdot \sigma$  and the identity permutation  $\iota$ , where  $\tau^{-1}$  denotes the inverse of  $\tau$ .<sup>7</sup> Hence, we can take as our input the permutation  $\pi = \tau^{-1} \sigma$ , and compute its distance from  $\iota$ . We call this formulation of the problem, *sorting by reversals*. Note also that any algorithm for the reversal distance between two strings that does not exploit a bounded-size alphabet must, as a special case, solve the reversal distance problem on permutations.

From an algebraic point of view, reversals generate the group of permutations under composition. Given an arbitrary group element  $\pi$ , we seek a shortest product of generators  $\rho_1 \rho_2 \cdots \rho_d$  that equals  $\pi$ .

*1.1. Related Work.* Little is known about reversal distance: even its computational complexity is open. The only reference to an algorithm appears to be in Watterson *et al.* [24], which gives the first definition of the problem, and a heuristic for computing reversal distance that is described in Section 2. Since there are so few references, we can give a fairly exhaustive coverage of related work.

<sup>4</sup> This notation is shorthand for  $\rho(i) = j, \rho(i+1) = j-1$ , etc. Outside interval  $[i, j]$ ,  $\rho$  leaves the elements unchanged.

<sup>5</sup> The *composition* of permutations  $\sigma$  and  $\rho$ , indicated by  $\sigma \cdot \rho$ , is a permutation  $\pi$  where  $\pi(i) = \sigma(\rho(i))$ .

<sup>6</sup> We also informally refer to  $d$  as the *inversion distance*.

<sup>7</sup> The *identity* permutation  $\iota$  is  $(1 \ 2 \ \cdots \ n)$ . The *inverse* of permutation  $\pi$  is the permutation  $\pi^{-1}$  satisfying  $\pi^{-1} \cdot \pi = \iota$ .

From the perspective of edit distance, the work of Wagner [23] is interesting. Wagner considers the problem of computing the minimum number of insertions, deletions, substitutions, and transpositions of adjacent characters, to convert one string to another, and shows that if the operations are restricted to deletion and transposition, the problem is NP-complete. If we restrict our problem to reversals of length two, in other words the adjacent transpositions of Wagner, the reversal distance between  $\sigma$  and  $\tau$  reduces to the so-called inversion number of  $\pi = \tau^{-1}\sigma$ , the number of pairs of  $i < j$  such that  $\pi_i > \pi_j$ , which is clearly computable in polynomial time (see [13, p. 11]). Tichy [22] also considers a variation of edit distance, but it is less closely related to our work. More recently, Schöniger and Waterman [20] present a heuristic for computing edit distance when only nonoverlapping inversions are allowed.

From the perspective of sorting, related work is by Gates and Papadimitriou [9]. They consider the problem of sorting a permutation by *prefix reversals*,<sup>8</sup> which are reversals of the form  $[1, i]$ , and derive bounds on the diameter of the problem. The *diameter* of the prefix reversal problem, which we denote by  $d_{\text{prefix}}(n)$ , is the maximum of the minimum number of prefix reversals to sort any permutation on  $n$  elements. Gates and Papadimitriou show that  $d_{\text{prefix}}(n) \leq \frac{5}{3}n + \frac{5}{3}$ , and that for infinitely many  $n$ ,  $d_{\text{prefix}}(n) \geq \frac{17}{16}n$ . Under the requirement that each element is reversed an even number of times, which may be appropriate if elements have an orientation (see Section 5), they show  $\frac{3}{2}n - 1 \leq d_{\text{prefix}}(n) \leq 2n + 3$ . In other work, Aigner and West [1] consider the diameter of sorting when the operation is reinsertion of the first element, and Amato *et al.* [2] consider a variation inspired by the problem of reversing trains on a track.

For our problem of sorting by unrestricted reversals, it appears tighter bounds on the diameter are possible. The heuristic of Watterson *et al.* [24] sorts any  $n$ -element permutation in  $n - 1$  reversals, so, writing  $d(n)$  for the diameter of our problem, we know  $d(n) \leq n - 1$ . From the other direction, Golan [10] has conjectured<sup>9</sup> that a particular  $n$ -element permutation, which we denote by  $\gamma^{(n)}$ , requires  $n - 1$  reversals, and has verified this for  $n$  up to 12. Recursively,

$$\gamma^{(n+1)} = \begin{cases} (1), & n \text{ is zero,} \\ (\gamma_1^{(n)} \ \gamma_2^{(n)} \ \cdots \ \gamma_{n-1}^{(n)} \ n+1 \ \gamma_n^{(n)}), & n \text{ is odd,} \\ (\gamma_1^{(n)} \ \gamma_2^{(n)} \ \cdots \ \gamma_{n-2}^{(n)} \ n+1 \ \gamma_n^{(n)} \ \gamma_{n-1}^{(n)}), & n \text{ is even.} \end{cases}$$

Using the lower bound developed in Section 3.2, we have verified the conjecture for  $n$  up to 200 when  $n \bmod 3 = 1$ ; for  $n \bmod 3 \neq 1$ , our computation showed  $\gamma^{(n)}$  requires  $n - 2$  reversals. If Golan's conjecture is true,  $d(n) = n - 1$ .

Note that studying the diameter of the problem, and algorithms that meet the

<sup>8</sup> This is also known as the *pancake flipping problem*.

<sup>9</sup> Golan's full conjecture is somewhat stronger: that, for every  $n$ ,  $\gamma^{(n)}$  and its inverse are the only permutations requiring  $n - 1$  reversals.

diameter, does not give a guarantee of quality of approximation. For example,  $\frac{1}{2}n + \frac{1}{2}$  is a lower bound on the diameter of our problem, and as we have indicated, the algorithm of Watterson *et al.* [24] uses no more than  $n - 1$  reversals; nevertheless, as we will show, there are permutations for which this algorithm performs arbitrarily poorly in ratio. Section 2 presents an approximation algorithm that does achieve a performance ratio of 2 for unrestricted reversals.

Finally, from the perspective of group theory, [6] and [11] are interesting. Even and Goldreich [6] show that, given a set of generators for a permutation group  $G$ , and a permutation  $\pi$ , determining the shortest product of generators that equals  $\pi$  is NP-hard.<sup>10</sup> Their reduction implies that the problem remains NP-hard even when every generator is its own inverse, as is the case in our problem. Jerrum [11] established that the problem is PSPACE-complete, and remains so when restricted to two generators.<sup>11</sup>

In our problem, the generator set is fixed. Thus, while these complexity results give us a sense of the problem, they do not imply the intractability of sorting by reversals. Nevertheless, we believe sorting by reversals is NP-complete. Section 5 indicates one possible direction for a proof.

*1.2. Overview.* In the next section, we present an approximation algorithm for sorting by reversals. We show that it never exceeds the minimum number by more than a factor of 2, and has a simple quadratic-time implementation.

Section 3 develops an exact algorithm using the branch-and-bound technique. The lower bound uses a relaxation to maximum-weight matchings, and linear programming.

Section 4 presents results from experiments with these algorithms. We study their performance on random permutations, and permutations generated by a fixed number of random reversals.

Section 5 concludes with some open problems and conjectures.

**2. An Approximation Algorithm.** Perhaps the most natural algorithm for sorting by reversals, suggested by Watterson *et al.* [24], is to bring element 1 into place, then element 2, and so on up to element  $n$ . Formally, at step  $i$ , perform reversal  $[i, \pi_i^{-1}]$ , if  $\pi_i \neq i$ . Once step  $n - 1$  is completed, element  $n$  must be in position  $n$ , so this sorts any  $n$ -element permutation in at most  $n - 1$  reversals.

While it is likely that permutations exist for every  $n$  that require  $n - 1$  reversals [10], which, if true, means this algorithm is worst-case optimal, for specific instances the algorithm can perform arbitrarily poorly. Consider, for example, the

<sup>10</sup> Even and Goldreich also show that computing the diameter of  $G$  is NP-hard. Determining whether there is a product equal to  $\pi$  is solvable in polynomial time [7].

<sup>11</sup> This result is best possible, since the case of a single generator can be solved in polynomial time. Jerrum also shows that the problem is polynomial-time solvable for any of the standard sets of generators for the symmetric and alternating groups: all transpositions or 2-cycles, all adjacent transpositions, all transpositions adjacent on the circle, all 3-cycles, and all 3-cycles with all pairs of disjoint transpositions.

permutation  $(n \ 1 \ 2 \ \dots \ n - 1)$ . Bringing 1 into place, then 2, and so on, uses  $n - 1$  reversals, yet the permutation can be sorted in two steps: reverse  $[1, n]$ , then  $[1, n - 1]$ . Thus, this algorithm can produce a solution  $\frac{1}{2}(n - 1)$  times longer than the shortest solution, for arbitrarily large  $n$ . Using the idea of a breakpoint, also introduced in [24], we show there is a simple algorithm guaranteed to use no more than twice the minimum number of reversals. To the best of our knowledge, this is the first constant-factor approximation for sorting by reversals.

In order to describe the algorithm, we first define some terminology.

A *breakpoint* of a permutation  $\pi$  is a pair of adjacent positions  $(i, i + 1)$  such that  $|\pi_{i+1} - \pi_i| \neq 1$ . In other words,  $(i, i + 1)$  forms a breakpoint if values  $\pi_i$  and  $\pi_{i+1}$  are not consecutively increasing or decreasing. To handle the boundaries, we let  $\pi_0$  have the value 0,  $\pi_{n+1}$  have the value  $n + 1$ , and allow  $i$  to range from 0 to  $n$  in the definition. Thus,  $(0, 1)$  is a breakpoint if  $\pi_1 \neq 1$ , and  $(n, n + 1)$  is a breakpoint if  $\pi_n \neq n$ . Notice that the identity permutation has no breakpoints, any other permutation has some breakpoint, and the number of breakpoints is at most  $n + 1$ .

When  $|\pi_{i+1} - \pi_i| = 1$ , we say values  $\pi_{i+1}$  and  $\pi_i$  are *adjacent*, and write  $\pi_{i+1} \sim \pi_i$ .

A *strip* of  $\pi$  is an interval  $[i, j]$  such that  $(i - 1, i)$  and  $(j, j + 1)$  are breakpoints, and no breakpoint lies between them. In other words, a strip is a maximal run of increasing or decreasing elements.

A reversal  $\rho$  affects the breakpoints of  $\pi$  only at the endpoints of  $\rho$ . (In the interior,  $\rho$  only makes an increasing pair  $(\pi_i, \pi_{i+1})$  decreasing, and vice versa.) Let us write  $\Phi(\pi)$  for the number of breakpoints in  $\pi$ , and, for a given reversal  $\rho$ , let

$$\Delta\Phi(\pi) = \Phi(\pi) - \Phi(\pi \cdot \rho).$$

Since a reversal  $[i, j]$  changes the adjacency of only two points, namely  $(i - 1, i)$  and  $(j, j + 1)$ , the only values  $\Delta\Phi(\pi)$  can take on are between  $-2$  and  $2$ . Since a solution must decrease the number of breakpoints from  $\Phi(\pi)$  to zero, a *greedy strategy* is to choose a reversal of maximum  $\Delta\Phi(\pi)$ , which achieves the greatest decrease. As any  $\pi \neq i$  has a reversal with  $\Delta\Phi(\pi) \geq 0$ , we can always achieve a decrease of 2, 1, or 0.

Figure 2 specifies our greedy algorithm. The algorithm removes zero breakpoints when there are no reversals that remove one or two, so it is not obvious that it terminates. With the rule “favor reversals that leave decreasing strips,” not only does the algorithm terminate, it exceeds the minimum by at most a factor of 2.

*2.1. Quality of the Approximation.* In the following, a strip  $[i, j]$  is *decreasing* if  $\pi_i, \pi_{i+1}, \dots, \pi_j$  is decreasing. We consider a strip of one element to be decreasing, except for  $\pi_0$  and  $\pi_{n+1}$ , which are always increasing. Thus, the identity permutation forms one increasing strip, extending from 0 to  $n + 1$ .

**LEMMA 1.** *Every permutation with a decreasing strip has a reversal that removes a breakpoint.*

```

algorithm GREEDY( $\pi$ ) begin
   $i := 0$ 
  while  $\pi$  contains a breakpoint do begin
     $i := i + 1$ 
    Let  $\rho_i$  be a reversal that removes the most breakpoints of  $\pi$ , resolving ties
      among those that remove one breakpoint in favor of reversals that leave a
      decreasing strip.
     $\pi := \pi \cdot \rho_i$ 
  end
  return  $i, (\rho_1, \rho_2, \dots, \rho_i)$ 
end
  
```

Fig. 2. The greedy algorithm.

PROOF. Consider the decreasing strip of  $\pi$  whose last element,  $\pi_i$ , is smallest. Element  $\pi_i - 1$  must be in an increasing strip (else  $\pi_i$  is not smallest), which lies either to the left or to the right of the strip containing  $\pi_i$ , as shown in Figure 3. In either case, the indicated reversal removes at least one breakpoint.  $\square$

LEMMA 2. Let  $\pi$  be a permutation with a decreasing strip. If every reversal that removes a breakpoint of  $\pi$  leaves a permutation with no decreasing strips,  $\pi$  has a reversal that removes two breakpoints.

PROOF. Again consider the decreasing strip of  $\pi$  containing the smallest element  $\pi_i$ . Case (b) of Figure 3 cannot occur, since  $\rho$  is a reversal that removes a breakpoint and leaves a decreasing strip. Thus, the increasing strip containing  $\pi_i - 1$  must be to the left of the strip containing  $\pi_i$ , as in case (a). Call the reversal of case (a),  $\rho_i$ .

Consider the decreasing strip of  $\pi$  whose first element,  $\pi_j$ , is greatest. Element  $\pi_j + 1$  must be in an increasing strip (else  $\pi_j$  is not greatest) that is to the right of the strip containing  $\pi_j$ , as otherwise, a reversal analogous to case (b) removes a

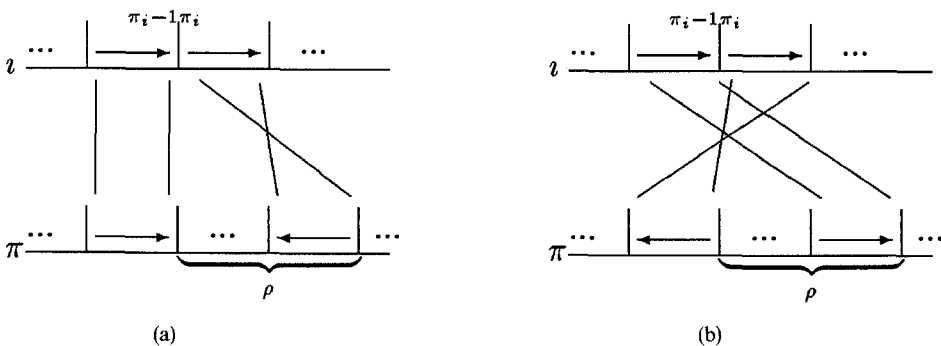


Fig. 3. A permutation  $\pi$  with a decreasing strip has a reversal  $\rho$  that removes a breakpoint. Element  $\pi_i$  is the smallest element that is in a decreasing strip.

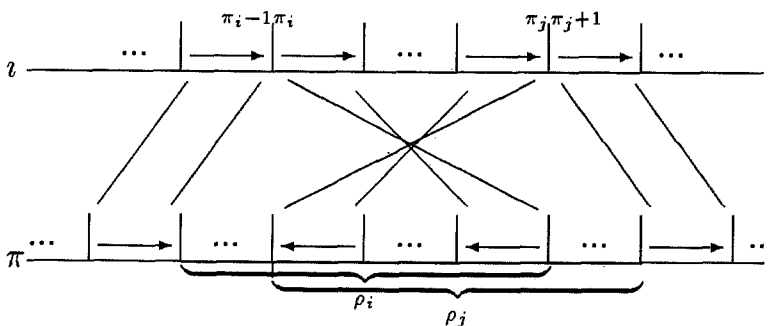


Fig. 4. If every reversal that removes a breakpoint of  $\pi$  leaves a permutation with no decreasing strips,  $\rho_i$  and  $\rho_j$  must overlap. Elements  $\pi_i$  and  $\pi_j$  are the smallest and largest in decreasing strips.

breakpoint and leaves a decreasing strip. Call  $\rho_j$  the reversal for  $\pi_j$  that is analogous to  $\rho_i$ .

Notice that  $\pi_j$  must lie in interval  $\rho_i$  and  $\pi_j + 1$  must lie outside, else  $\rho_i$  leaves a decreasing strip. Similarly,  $\pi_i$  must lie in  $\rho_j$  and  $\pi_i - 1$  outside, else  $\rho_j$  leaves a decreasing strip. The situation is as shown in Figure 4. Intervals  $\rho_i$  and  $\rho_j$  overlap.

We now argue that not only do  $\rho_i$  and  $\rho_j$  overlap, they must be the same interval. For suppose  $\rho_i - \rho_j$  is not empty. If it contains a decreasing strip, reversal  $\rho_j$  leaves a decreasing strip, and if it contains an increasing strip, reversal  $\rho_i$  leaves a decreasing strip. Similarly, interval  $\rho_j - \rho_i$  must be empty, which implies  $\rho_i = \rho_j$ .

Since reversal  $\rho_i$  removes the breakpoint on its left, and reversal  $\rho_j$  removes the breakpoint on its right, and as these breakpoints are distinct, reversal  $\rho = \rho_i = \rho_j$  removes two breakpoints. □

LEMMA 3. *The greedy algorithm sorts a permutation  $\pi$  with a decreasing strip in at most  $\Phi(\pi) - 1$  reversals.*

PROOF. The proof is by induction on  $\Phi(\pi)$ . If  $\pi$  has a decreasing strip,  $\Phi(\pi) \geq 2$ . When  $\Phi(\pi) = 2$ ,  $\pi$  has a unique reversal  $\rho$  that removes the two breakpoints and sorts  $\pi$ . Since GREEDY will choose  $\rho$ , it sorts  $\pi$  in one reversal, and the basis holds.

Suppose the lemma holds for all  $\pi'$  with less than  $\Phi(\pi)$  breakpoints. Since  $\pi$  has a decreasing strip, by Lemma 1 there is a reversal  $\rho$  that removes at least one breakpoint of  $\pi$ . Thus, the first step of GREEDY will transform  $\pi$  into a permutation  $\pi'$  with at most  $\Phi(\pi) - 1$  breakpoints. If  $\pi'$  has a decreasing strip, GREEDY sorts it in at most  $\Phi(\pi) - 2$  reversals by the induction hypothesis, which sorts  $\pi$  in at most  $\Phi(\pi) - 1$  reversals.

Now consider a  $\pi'$  with no decreasing strips. We argue that  $\Phi(\pi') = \Phi(\pi) - 2$ . For suppose  $\Phi(\pi') = \Phi(\pi) - 1$ , the only other possibility. Since GREEDY chooses a reversal that removes the most breakpoints, every reversal that removes a



breakpoint must remove exactly one breakpoint. Since, in such an event, GREEDY chooses a reversal that leaves a decreasing strip whenever possible, every available reversal that removes a breakpoint must leave no decreasing strips. However, by Lemma 2, this implies  $\pi$  has a reversal that removes two breakpoints, a contradiction. Thus  $\Phi(\pi') = \Phi(\pi) - 2$ .

Every reversal on a permutation with no decreasing strips creates a decreasing strip, which implies GREEDY will transform  $\pi'$  to a permutation  $\pi''$  with a decreasing strip. Moreover,  $\Phi(\pi'') \leq \Phi(\pi') = \Phi(\pi) - 2$ . By induction, GREEDY sorts  $\pi''$  in at most  $\Phi(\pi) - 3$  reversals. Since GREEDY transformed  $\pi$  to  $\pi''$  in two steps, this sorts  $\pi$  in at most  $\Phi(\pi) - 1$  reversals. □

**THEOREM 1.** *The greedy algorithm sorts every permutation  $\pi$  in at most  $\Phi(\pi)$  reversals.*

**PROOF.** If  $\pi$  has a decreasing strip, by Lemma 3, GREEDY sorts it within  $\Phi(\pi)$  reversals. If  $\pi$  has no decreasing strip, any reversal chosen by GREEDY transforms  $\pi$  to a permutation  $\pi'$  with a decreasing strip such that  $\Phi(\pi') \leq \Phi(\pi)$ . By Lemma 3, GREEDY sorts  $\pi'$  in at most  $\Phi(\pi') - 1$  reversals, which sorts  $\pi$  in at most  $\Phi(\pi)$  reversals. □

Since  $\Phi(\pi) \leq n + 1$ , Theorem 1 implies that the greedy algorithm terminates in  $O(n)$  iterations. In the next section, we consider how to implement an iteration; here we simply note that an iteration runs in polynomial time.

An algorithm for an optimization problem that runs in polynomial time and delivers a solution whose value is within a factor  $\alpha$  of optimal is known as an  $\alpha$ -approximation algorithm. An immediate consequence of Theorem 1 is the following.

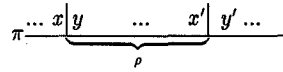
**COROLLARY 1.** *The greedy algorithm is a 2-approximation algorithm for sorting by reversals.*

**PROOF.** Write  $\text{OPT}(\pi)$  for the minimum number of reversals to sort a permutation  $\pi$ , and  $\text{GREEDY}(\pi)$  for the number taken by the greedy algorithm. Since a solution must remove all breakpoints, and any reversal can remove at most two,

$$\text{OPT}(\pi) \geq \left\lceil \frac{1}{2} \Phi(\pi) \right\rceil \geq \frac{1}{2} \text{GREEDY}(\pi). \quad \square$$

We do not know whether the bound of Corollary 1 is tight.

**2.2. Time and Space.** How much time does an iteration of the greedy algorithm take? As there are in general  $\binom{n}{2}$  reversals to consider, a naive implementation could take  $O(n^2)$  time per iteration, or  $O(n^3)$  time in total. By considering the form of reversals that remove breakpoints, we can find the greedy reversal for an iteration in  $O(n)$  time, which yields an  $O(n^2)$ -time algorithm.



**Fig. 5.** The form of a reversal that removes two breakpoints. Vertical bars denote breakpoints,  $x \sim x'$ , and  $y \sim y'$ .

A reversal that removes two breakpoints must have each endpoint at a breakpoint, and must create two adjacencies. Let us denote the left endpoint by  $(i, i + 1)$  and the right endpoint by  $(j, j + 1)$ . Then reversal  $[i + 1, j]$  removes two breakpoints iff  $(i, i + 1)$  and  $(j, j + 1)$  are breakpoints,  $\pi_i \sim \pi_j$ , and  $\pi_{i+1} \sim \pi_{j+1}$ . This is illustrated in Figure 5.

We can search for a reversal of this form as follows. Scanning  $\pi$ , we identify each breakpoint  $(i, i + 1)$ . If this is the left end of such a reversal, there must be a position  $j > i$  such that  $(j, j + 1)$  is a breakpoint,  $\pi_j \sim \pi_i$ , and  $\pi_{j+1} \sim \pi_{i+1}$ . There are two possible values for  $\pi_j$ , namely  $\pi_i - 1$  and  $\pi_i + 1$ . Given  $\pi^{-1}$ , we can find the positions where  $\pi_i - 1$  and  $\pi_i + 1$  occur in  $O(1)$  time. If either position meets the criteria above, we have found a reversal that removes two breakpoints. As there are  $O(n)$  candidates for the left endpoint, and  $\pi^{-1}$  can be computed in  $O(n)$  time, this finds a reversal that removes two breakpoints (if one exists) in  $O(n)$  time.

If there is no reversal that removes two breakpoints, the greedy algorithm considers reversals that remove one breakpoint. A reversal that removes one breakpoint must have an endpoint at the breakpoint it removes; the other end may or may not lie at a breakpoint as shown in Figure 6. Given  $\pi^{-1}$ , we can find a reversal of form (a) or (b) in  $O(1)$  time per breakpoint, as explained above. The only question is how to determine efficiently whether the reversal leaves a decreasing strip.<sup>12</sup>

Consider a reversal of form (a). Reversal  $[i, j]$  leaves a decreasing strip iff

- $[1, i)$  contains a decreasing strip other than  $x$ , or
- $[i, j]$  contains an increasing strip<sup>13</sup> other than  $x'$ , or
- $(j, n]$  contains a decreasing strip, or
- $xx'$  is decreasing.

The only difficulty is in determining whether an interval contains an increasing strip, or a decreasing strip, in  $O(1)$  time. We can solve this by forming an array,



**Fig. 6.** The form of a reversal that removes one breakpoint. In the figure,  $x \sim x' \sim x'' \sim x'''$  and  $y \not\sim z$ . Mirror images of (a) and (b) are considered to be the same form.

<sup>12</sup> Recall that the greedy algorithm breaks ties among reversals that remove one breakpoint by favoring reversals that leave decreasing strips.

<sup>13</sup> A strip is *increasing* if its elements are strictly increasing, or it contains a single element. Thus, a single-element strip is both increasing and decreasing.

$down[i]$ , that gives the position of the left end of the leftmost decreasing strip beginning at, or to the right of, position  $i$ . Then interval  $[a, b]$  contains a decreasing strip iff  $down[a] \leq b$ . Similarly, we can form an array  $up[i]$  that gives the left end of the leftmost increasing strip at, or to the right of,  $i$ . Both arrays can be computed from  $\pi$  in  $O(n)$  time, for example by the recurrence,

$$down(i) = \begin{cases} n + 1, & i > n, \\ i, & i \leq n, \pi_i \not\prec \pi_{i-1}, \pi_i \not\prec \pi_{i+1}, \\ i, & i \leq n, \pi_i \not\prec \pi_{i-1}, \pi_i \sim \pi_{i+1}, \pi_i > \pi_{i+1}, \\ down(i + 1), & \text{otherwise.} \end{cases}$$

Thus, if there is a reversal of form (a), we can find it in  $O(n)$  time, whether or not we require that it leave a decreasing strip. We can also search for a reversal of form (b) in  $O(n)$  time, using the same technique.

If there is no reversal that removes a breakpoint, the greedy algorithm chooses a reversal that does not increase  $\Phi(\pi)$ . One such reversal is  $[i, \pi_i^{-1}]$ , where  $i$  is the smallest position such that  $\pi_i \neq i$ . Notice that this reversal cannot increase  $\Phi(\pi)$  since it always removes breakpoint  $(i - 1, i)$ .

To summarize, we find a greedy reversal as follows:

- (1) Compute  $\pi^{-1}$ ,  $down$ , and  $up$ .
- (2) Search for a reversal that removes two breakpoints.
- (3) If none exists, search for a reversal that removes one breakpoint *and* leaves a decreasing strip.
- (4) If none exists, search for a reversal that removes one breakpoint.
- (5) If none exists, bring the smallest out-of-place element into position.

Each step can be performed in  $O(n)$  time; which gives an  $O(n^2)$ -time implementation of the greedy algorithm (and with more care,  $O(n + \Phi^2(\pi))$  time can be achieved). We suspect that an  $O(n \log n)$ -time implementation may be possible; our experience, however, suggests that the approximation algorithm will be far from the dominant step in practice; as we discuss in Section 4.

**3. Exact Algorithm.** In the preceding section, we obtained an algorithm that comes close to the optimum by applying a greedy strategy: of all reversals, select one that removes the most breakpoints. To obtain an algorithm that reaches the optimum, we use a *branch-and-bound* strategy: consider all reversals, and eliminate those that cannot lead to an optimal solution.

Figure 7 shows the form of our branch-and-bound algorithm. We maintain three global variables: *bound*, a dynamic upper bound on the solution value; *best*, an array of reversals that sort the permutation in *bound* steps; and *current*, the series of reversals currently under consideration. At the start, we initialize *bound*

```

global bound, current[1..n], best[1..n]

algorithm BRANCHANDBOUND( $\pi$ ) begin
  bound, best := UPPERBOUND( $\pi$ )
  SEARCH( $\pi$ , 0)
  return bound, best
end

algorithm SEARCH( $\pi$ , depth) begin
  if  $\pi$  is the identity permutation then
    if depth < bound then bound, best := depth, current
  else
    for each reversal  $\rho$  in order of decreasing  $\Delta\Phi(\pi)$  do
      if LOWERBOUND( $\pi \cdot \rho$ ) + depth + 1 < bound then begin
        current[depth + 1] :=  $\rho$ 
        SEARCH( $\pi \cdot \rho$ , depth + 1)
      end
    end
  end

```

Fig. 7. The branch-and-bound algorithm.

and *best* to values obtained from an upper-bound algorithm.<sup>14</sup> The algorithm we use is essentially GREEDY with a fixed-depth look-ahead, and is described in Section 3.3.

After obtaining an upper bound, we explore a tree of subproblems depth-first. Each invocation of SEARCH corresponds to a node of the tree and is labeled with  $\pi$ , a permutation to be sorted, and *depth*, the number of edges from the root to the node. Array *current* is maintained as a stack by SEARCH, and holds the reversals on the path from the root to the current node. We chose a depth-first strategy for traversing the tree as this uses a polynomial amount of space, even when the tree is of exponential size, since space, not time, is often the limiting resource.

Examining all reversals yields a very large tree: with  $\binom{n}{2}$  children per node, and a height of  $n - 1$ , there are  $O(n^{2^n}/2^n)$  nodes. In Section 5, we state several conjectures on the form of a solution, which, if true, reduce the children per node from  $O(n^2)$  to  $O(\Phi^2(\pi))$ . Lacking a proof of these conjectures, the two means we have to reduce the size of the search trees are ordering children, and computing lower bounds. The algorithm of Figure 7 orders children by decreasing  $\Delta\Phi$ , on the assumption that the optimal solution uses reversals of greatest  $\Delta\Phi$ . By trying such reversals first, we hope to lower our upper bound quickly, to prune subtrees early on. We now explain how the lower bound is computed.

**3.1. A Lower Bound from Matchings.** As stated in the proof of Corollary 1, a simple lower bound on  $\text{OPT}(\pi)$  is  $\lceil \Phi(\pi)/2 \rceil$ . While this is sufficient to prove an approximation factor of 2, it is extremely weak. It assumes every

<sup>14</sup> We use  $x_1, \dots, x_n := e_1, \dots, e_n$  as shorthand for the parallel assignments  $x_i := e_i$ . Function UPPERBOUND, like GREEDY, returns two values: an integer, followed by a list of reversals.

breakpoint of  $\pi$  can be eliminated by a reversal that removes two breakpoints, which can rarely be achieved. To obtain a better bound, we ask, for a given permutation, how many breakpoints can possibly be eliminated by reversals that remove two breakpoints?

A pair of breakpoints  $p = (i, i + 1)$  and  $q = (j, j + 1)$ , with values  $(\pi_i, \pi_{i+1}) = (x, y)$  and  $(\pi_j, \pi_{j+1}) = (x', y')$ , can be eliminated in one reversal iff  $x \sim x'$  and  $y \sim y'$ . This holds whether  $p$  is to the left or to the right of  $q$ . The only requirement is that  $x$  and  $y$  occur in the same order as  $x'$  and  $y'$ .

Notice that such a reversal also affects other pairs of breakpoints. A pair with values  $(x, y)$  and  $(y', x')$ , which cannot be eliminated immediately because adjacent values are not in the same order, can be eliminated in one step if preceded by a reversal that contains exactly one breakpoint of the pair. Such a reversal transforms the pair to the preceding case.

In general, determining when a collection of  $2m$  breakpoints can be eliminated by a sequence of  $m$  reversals appears difficult. (In Section 5, we conjecture it is NP-complete.) To obtain a lower bound that can be efficiently computed, we ignore dynamic information about the order and interaction of reversals. The static information we retain is simply the adjacency of values between breakpoints, which can be represented by a graph.

Figure 8 shows the construction. Each breakpoint of  $\pi$  is mapped to a vertex of  $G(\pi)$ . We place an edge between breakpoints  $p$  and  $q$  if either of the above two cases apply. Effectively, if two breakpoints can be eliminated by one reversal, possibly after a sequence of other reversals that eliminate two breakpoints, they share an edge. Note that the order of the two values at a breakpoint is not important in the construction.

Since each edge models a reversal, and performing the reversal removes both endpoints, a series on reversals on  $\pi$  that each eliminate two breakpoints corresponds to a set of vertex-disjoint edges in  $G(\pi)$ . A set of vertex-disjoint edges is called a *matching*. The key property of  $G(\pi)$  is that the most reversals we can possibly perform on  $\pi$  that each remove a pair of breakpoints, without performing any intervening reversals that remove less than two breakpoints, is the size of a maximum-cardinality matching of  $G(\pi)$ .

Let  $m$  be the number of vertices in a maximum-cardinality matching of  $G(\pi)$ , in other words, twice the number of edges in the matching. How many reversals must be performed to remove the remaining  $\Phi(\pi) - m$  breakpoints of  $\pi$ ? The best we can do is to expend a reversal that removes one breakpoint to set up a reversal that removes two breakpoints. Notice that we cannot remove one breakpoint,

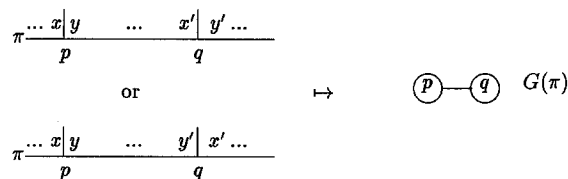


Fig. 8. Breakpoints  $\{x, y\}$  and  $\{x', y'\}$  share an edge in  $G(\pi)$  iff  $x \sim x'$  and  $y \sim y'$ .

then two, then two again. (A reversal that removes one breakpoint can affect only one additional breakpoint. This implies the third reversal must have been available from the start, which contradicts that the matching is of maximum cardinality.) In short, the best we can do is to remove three breakpoints in two reversals. This gives a lower bound of

$$(1) \quad \lceil \frac{1}{2}m + \frac{2}{3}(\Phi(\pi) - m) \rceil,$$

which has the extreme value  $\lceil \frac{2}{3}\Phi(\pi) \rceil$ .

We can construct  $G(\pi)$  from  $\pi$  in  $O(n)$  time. (Certainly the  $O(n)$  breakpoints of  $\pi$  can be determined in  $O(n)$  time. Moreover, every breakpoint is incident to a constant number of edges, since the only values adjacent to  $x$  are  $x + 1$  and  $x - 1$ . So, with the help of  $\pi^{-1}$ , we can determine all edges in  $O(n)$  time as well.) A maximum-cardinality matching of a graph with  $V$  vertices and  $E$  edges can be computed in  $O(E\sqrt{V})$  time [15]. Thus, since  $V$  and  $E$  for  $G(\pi)$  are both  $O(n)$ , we can evaluate the lower bound of (1) in  $O(n^{3/2})$  time.

*3.2. A Family of Lower Bounds.* We can improve the lower bound further, by considering 3-tuples of breakpoints, 4-tuples of breakpoints, and so on.

Let us call a reversal that eliminates  $k$  breakpoints, a  $k$ -move. Thus, a 2-move is a reversal that eliminates two breakpoints, and a  $(-2)$ -move is a reversal that creates two breakpoints.

In general, for  $k \geq 3$ , we define a  $k$ -move as follows. Over all permutations, consider all series of reversals that eliminate  $k$  breakpoints. A  $k$ -move, for  $k \geq 3$ , is a shortest series that eliminates a set of  $k$  breakpoints, given that no 2-, 3-, up to  $(k - 1)$ -moves are available on the set. For example, a 3-move is a 1-move followed by a 2-move. (Notice that this arose in the analysis of lower bound (1).)

The following lemma characterizes the structure of a  $k$ -move.

LEMMA 4. *For  $k \geq 3$ , a  $k$ -move is a series of  $k - 1$  reversals, that decomposes into either*

- (i) *a 1-move followed by a  $(k - 1)$ -move, or*
- (ii) *a 0-move followed by an  $i$ -move and a  $j$ -move, where  $i + j = k$ .*

PROOF. Any series of reversals begins with a 2-, 1-, 0-,  $(-1)$ -, or  $(-2)$ -move. By definition, a  $k$ -move for  $k \geq 3$  cannot begin with a 2-move. Furthermore, any series that creates breakpoints is not among the shortest. Thus, 1-moves and 0-moves are the only candidates for the first reversal in a  $k$ -move.

Consider a series that begins with a 1-move. The 1-move can change the values of at most two breakpoints, namely, those at its endpoints. One of the two breakpoints is eliminated by the 1-move. The other can at best be eliminated in a  $(k - 1)$ -move. Note that the  $k - 1$  breakpoints remaining cannot be eliminated by two or more higher-order moves, as any second higher-order move would be available from the start, contradicting the definition of a  $k$ -move.

Now consider a series that begins with a 0-move. This move can again affect

the values of at most two breakpoints, thereby setting up at most two higher-order moves. The  $k$  breakpoints again cannot be eliminated by three or more higher-order moves, after performing the initial 0-move, since a third higher-order move would be available initially. That a  $k$ -move uses  $k - 1$  reversals follows by induction.  $\square$

With this decomposition, we can characterize the breakpoints in a  $k$ -move.

LEMMA 5. *For  $k \geq 2$ , the pairs of values at the breakpoints eliminated by a  $k$ -move have the form*

$$(2) \quad \{x_1, x'_2\} \{x_2, x'_3\} \cdots \{x_{k-1}, x'_k\} \{x_k, x'_1\},$$

where  $x_i \sim x'_i$  for  $1 \leq i \leq k$ .

PROOF. Notice that the lemma holds for  $k = 2$ , which corresponds to the picture of Figure 8. Assume the lemma then for all  $k' \leq k$ . We show it holds for  $k + 1$ .

By Lemma 4, a  $(k + 1)$ -move decomposes into a 1-move followed by a  $k$ -move, or a 0-move followed by an  $i$ -move and a  $j$ -move. Consider the case of an initial 1-move.

This move eliminates one of the  $k + 1$  breakpoints, and brings the remaining  $k$  breakpoints into the configuration of a  $k$ -move. By induction, the values at these  $k$  breakpoints have the form of (2). Notice that this form is unchanged by a rotation of the breakpoints, i.e., a renaming of the form  $x_i \mapsto x_{i \oplus d}$  and  $x'_i \mapsto x'_{i \oplus d}$ , for any  $d$ , where  $i \oplus d$  denotes  $((i + d) \bmod k) + 1$ . Thus, we may assume without loss of generality, that the 1-move affects the values in (2) by bringing  $x_1, x'_2$  together into a breakpoint, and some other pair of values  $y, y'$  together to create an adjacency  $y \sim y'$ . This  $(k + 1)$ -move then has the form

$$\{x_1, y'\} \{y, x'_2\} \{x_2, x'_3\} \cdots \{x_k, x'_1\},$$

which is the same form as (2).

Now consider the case of a 0-move that sets up an  $i$ -move and a  $j$ -move. By induction, the  $i$ -move has the form

$$\{x_1, x'_2\} \{x_2, x'_3\} \cdots \{x_i, x'_1\},$$

and the  $j$ -move has the form

$$\{y_1, y'_2\} \{y_2, y'_3\} \cdots \{y_j, y'_1\}.$$

Without loss of generality, assume the 0-move brings  $x_1, x'_2$  together and  $y_1, y'_2$

together. Since the form of (2) is unchanged by renaming  $x_i \mapsto x'_i$  and  $x'_i \mapsto x_i$ , we may further assume that the 0-move brings these values together by touching breakpoints

$$\{x_1, y_1\} \{x'_2, y_2\}.$$

The  $(k + 1)$ -move then has the form

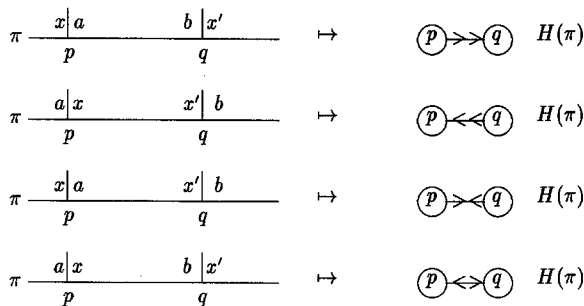
$$\{x_2, x'_3\} \{x_3, x'_4\} \cdots \{x_i, x'_1\} \{x_1, y_1\} \{y_j, y'_1\} \{y_{j-1}, y'_j\} \cdots \{y_2, y'_3\} \{x'_2, y_2\}.$$

As the reader may verify, this is the same form as (2). □

We now describe how to construct a graph  $H(\pi)$  that allows us to efficiently identify sets of breakpoints of form (2). In the construction, breakpoints of  $\pi$  are mapped to vertices of  $H(\pi)$ , and pairs of breakpoints that share an adjacency, such as  $\{x, a\} \{x', b\}$  where  $x \sim x'$ , are mapped to edges, as shown in Figure 9.

Edges of  $H(\pi)$  are directed, but not in the standard sense. An edge touching  $v$  and  $w$  contributes to either the in- or out-degree of  $v$  and  $w$ . In a directed graph, an edge  $(v, w)$  that contributes to the out-degree of  $v$  necessarily contributes to the in-degree of  $w$ , and vice versa. However, in what we call a *bidirected graph*, there are two more possibilities:  $(v, w)$  may contribute to the in-degree of both  $v$  and  $w$ , or to the out-degree of both  $v$  and  $w$ .

This gives rise to the four types of edges of Figure 9. We indicate the direction of an edge with double-ended arrows. When drawing an edge incident to  $v$ , we place an arrowhead pointing into  $v$ , at the end touching  $v$ , if the edge contributes to  $v$ 's in-degree. Otherwise, we direct the arrowhead out of  $v$ .



**Fig. 9.** Construction of graph  $H(\pi)$ . Values  $x$  and  $x'$  are adjacent, and induce an edge that contributes to the in- or out-degree of  $p$ , and the in- or out-degree of  $q$ , depending on whether the value is to the right or to the left of the breakpoint.



The utility of this construction is in the correspondence between cycles in the graph and  $k$ -moves on the permutation, as summarized in the following lemma. A  $k$ -cycle in a bidirected graph is a series of edges

$$(v_1, v_2) (v_2, v_3) \cdots (v_k, v_1)$$

such that the  $v_i$  are distinct, and every  $v_i$  has in- and out-degree 1.

LEMMA 6. *The sets of breakpoints of  $\pi$  whose values have the form of  $k$ -moves are in one-to-one correspondence with the  $k$ -cycles of  $H(\pi)$ .*

PROOF. By Lemma 5, the values in a  $k$ -move have the form

$$\{x_1, x'_2\} \{x_2, x'_3\} \cdots \{x_k, x'_1\}.$$

In the  $i$ th breakpoint,  $\{x_i, x'_{i \oplus 1}\}$ , value  $x'_{i \oplus 1}$  is adjacent to value  $x_{i \oplus 1}$  of the  $(i \oplus 1)$ th breakpoint. Whether  $x'_i$  is to the left or right of  $x'_{i \oplus 1}$  in the  $i$ th breakpoint, breakpoint  $i$  is linked in  $H(\pi)$  to breakpoints  $i \oplus 1$  and  $i \ominus 1$  by edges that contribute exactly once to its in- and out-degree.

Similarly, any  $k$ -cycle of  $H(\pi)$  describes a set of  $k$  breakpoints with the property that every breakpoint in the set has values that are adjacent to the preceding and succeeding breakpoints on the cycle. As every vertex of the cycle has in- and out-degree 1, these values, by the construction, are distinct. By Lemma 5, this is the form of a  $k$ -move.  $\square$

We now have the tools to generalize the lower bound of Section 3.1. In outline, we construct a hypergraph  $G^{(k)}(\pi)$  whose vertices correspond to breakpoints, but whose edges are sets of up to  $k$  vertices that correspond to  $k'$ -moves for  $k' \leq k$ . A series of moves on  $\pi$  maps to a matching of  $G^{(k)}$ , where a matching of a hypergraph is a collection of vertex-disjoint edges. Choosing a  $k$ -move corresponds to performing a series of  $k - 1$  reversals. We weight edges by the number of reversals they represent, and seek, as before, a maximum-weight matching. However, computing a maximum-weight matching of a hypergraph is in general NP-complete [8]. We express the matching problem as an integer programming problem, and relax the integrality constraint to obtain a linear programming problem. This gives a somewhat weaker lower bound that is computable in polynomial time.

This approach is summarized in the following theorem. The  $k$ -girth of a graph is the length of a shortest cycle of more than  $k$  edges. If the graph does not contain such a cycle, we define its  $k$ -girth to be  $n + 1$ , where  $n$  is the number of vertices.

**THEOREM 2.** *Let  $(V, E)$  be the graph  $G^{(k)}(\pi)$ , let  $g$  be the  $k$ -girth of  $H(\pi)$ , and let  $\mathcal{L}_k(\pi)$  be the solution value of the linear program*

$$\begin{aligned} & \text{minimize} && \frac{g-1}{g} \Phi(\pi) - \sum_{e \in E} \frac{g-|e|}{g} x_e, \\ & \text{subject to} && 0 \leq x_e \leq 1, \quad \text{for all } e \in E, \\ & && \sum_{e|v \in e} x_e \leq 1, \quad \text{for all } v \in V. \end{aligned}$$

Then  $\text{OPT}(\pi) \geq \lceil \mathcal{L}_k(\pi) \rceil$ .

**PROOF.** As in the lower bound of Section 3.1, the only characteristic of a reversal that we consider are the values at its endpoints. This means we ignore the effect of a reversal on its interior, namely, that it changes the relative order of elements.

This being the case, we first argue that to demonstrate our lower bound, the only series that we have to consider are those that do not create breakpoints. For suppose a reversal in a series creates a breakpoint. Eventually this breakpoint must be removed, and the best we can possibly do is to eliminate it with a 2-move, which allows us to remove one more breakpoint. The best we can then have achieved is to remove one breakpoint in two reversals. This is worse than any of the higher-order moves we consider in the lower bound, which always remove  $k$  breakpoints in  $k - 1$  reversals. Admittedly, the  $(-1)$ -move and subsequent 2-move, by changing the relative order of elements, may have made some advantageous moves possible in an actual series, but we have already accounted for such effects by ignoring the order of values at breakpoints in our graph representation.

Thus, it suffices to consider series that do not create breakpoints. Such a series operates only on breakpoints in the original permutation, by moving values from one breakpoint to another, so as to create adjacencies. Decomposing the series into higher-order moves, every  $k'$ -move, where  $k' \leq k$ , maps to an edge of  $G^{(k)}$ , no matter where it occurs in the series. Moreover, the edges so identified in  $G^{(k)}$  are vertex disjoint, and form a matching  $M$ . Thus, the number of reversals taken by  $k'$ -moves in the series, where  $k' \leq k$ , is

$$\sum_{e \in M} (|e| - 1).$$

The remaining moves of the series are  $k'$ -moves where  $k' > k$ . By Lemma 6, every one of these moves maps to a cycle of  $H(\pi)$ , so the smallest  $k' > k$  for which the series contains a  $k'$ -move is at least  $g$ , the  $k$ -girth of  $H(\pi)$ . Thus, the number of reversals taken by  $k'$ -moves in the series, where  $k' > k$ , is at least

$$(g-1) \left\lceil \frac{\Phi(\pi) - \sum_{e \in M} |e|}{g} \right\rceil + \left( \Phi(\pi) - \sum_{e \in M} |e| \right) \pmod{g},$$

which is

$$\left\lceil \frac{g-1}{g} \left( \Phi(\pi) - \sum_{e \in M} |e| \right) \right\rceil.$$

Thus, the total number of reversals in an optimal series is at least

$$\min_M \left[ \sum_{e \in M} (|e| - 1) + \frac{g-1}{g} \left( \Phi(\pi) - \sum_{e \in M} |e| \right) \right],$$

which is

$$\left\lceil \min_M \left\{ \frac{g-1}{g} \Phi(\pi) - \sum_{e \in M} \frac{g-|e|}{g} \right\} \right\rceil.$$

This is equivalent to finding a matching in  $G^{(k)}$  of maximum total weight, where the weight of an edge  $e$  is  $(g - |e|)/g$ . We can express this as an integer programming problem. For each edge  $e$  of  $G^{(k)}$ , we have a variable  $x_e$ , that takes on the values 0 or 1. Selecting  $e$  is encoded by assigning  $x_e$  the value 1. We can ensure that the assignment represents a matching by requiring

$$\sum_{e|v \in e} x_e \leq 1,$$

for every vertex  $v$ . Extending the domain of  $x_e$  to real values between 0 and 1 results in the linear programming problem of the theorem. □

Notice that the lower bound of Theorem 2 has the extreme value  $\lceil ((g-1)/g)\Phi(\pi) \rceil$ . When the  $k$ -girth  $g$  is large (as is the case with the Golan permutation  $\gamma^{(n)}$  even for small  $k$ ), this can be as great as  $\Phi(\pi)$ , which meets the upper bound of Theorem 1.

How much time does it take to evaluate  $\mathcal{L}_k$ ? There are three tasks:

- (1) Constructing  $H$  and computing its  $k$ -girth.
- (2) Constructing  $G^{(k)}$  and its associated linear program.
- (3) Solving the linear program.

Constructing  $H$  takes time  $O(n)$ . There are  $O(n)$  breakpoints, and each breakpoint has at most four in-edges and out-edges, which can be identified in  $O(1)$  time using  $\pi^{-1}$ .

We can compute the  $k$ -girth of  $H$  in  $O(4^k n^2)$  time, as follows. A shortest cycle of more than  $k$  edges, that contains a fixed vertex  $v$ , is a path  $P$  of  $k$  edges from  $v$  to some vertex  $w$ , followed by a shortest path from  $w$  to  $v$  that does not visit any other vertices on  $P$ . Paths in a bidirected graph such as  $H$  alternate in- and out-edges: if we enter a vertex by an in-edge, we must leave by an out-edge, and vice versa. As every vertex of  $H$  has in- and out-degree at most four, there are at

most  $4^k$  paths of length  $k$  from a fixed vertex  $v$ . For each path  $P$ , we can mark its vertices, and compute a shortest return path from its end  $w$  back to  $v$ , taking care not to visit marked vertices. This shortest path can be found by a breadth-first search from  $w$  in  $O(n)$  time. Repeating for all start vertices  $v$ , all paths  $P$ , and recording the minimum over all cycles found, takes time  $O(4^k n^2)$ .

Similarly, we can construct the edges of  $G^{(k)}$  in  $O(4^k n)$  time, by enumerating the cycles of  $H$  of  $k$  or fewer edges in a depth-first search.<sup>15</sup> Space for all edges is  $O(k4^k n)$ .

The resulting linear program has  $O(4^k n)$  variables, and  $O(n)$  constraints. Writing  $L(a, b)$  for the time to solve a linear program of  $a$  variables and  $b$  constraints, the linear programming problem takes  $O(L(4^k n, n))$  time. This dominates the time to compute the lower bound. Thus, for any fixed  $k$ ,  $\mathcal{L}_k$  can be computed in  $O(L(n, n))$  time.

**3.3. A Family of Upper Bounds.** As well as a lower bound on the solution value, our exact algorithm requires an upper bound. The simplest approach is to use the approximation algorithm of Section 2, but for large  $n$  this gives too weak a bound to prune away much of the search tree.

Consider a series of  $k$  reversals that removes the most breakpoints among series of that length. The greedy strategy of the approximation algorithm is really based on the observation that, once we are  $k$  reversals away from sorting a permutation, such a series is optimal. GREEDY corresponds to the case  $k = 1$ .

Such a series can be found by looking ahead  $k$  reversals, and this search can be made tractable by again employing branch-and-bound. The basic form of the computation is identical to BRANCHANDBOUND, except that the recursion is stopped at depth  $k$ . The two requirements are a lower bound on the number of breakpoints that can be eliminated in a series of  $k$  reversals, and a method for computing an upper bound on the number eliminated by an optimal extension of a partial series.

We can compute the lower bound by running GREEDY. Computing an upper bound is a little more difficult, but can be tackled by the methods of the previous section, as summarized in the following theorem.

**THEOREM 3.** *Let  $(V, E)$  be the graph  $G^{(k+1)}(\pi)$ , and let  $\mathcal{U}_k(\pi)$  be the solution value of the linear program*

$$\begin{aligned}
 & \text{maximize} && k + \sum_{e \in E} x_e, \\
 & \text{subject to} && 0 \leq x_e \leq 1, && \text{for all } e \in E, \\
 & && \sum_{e|v \in e} x_e \leq 1, && \text{for all } v \in V, \\
 & && \sum_{e \in E} (|e| - 1)x_e \leq k.
 \end{aligned}$$

<sup>15</sup> Note that determining the edges of  $G^{(k)}$  by examining all  $k$ -sets would take  $O(n^k)$  time.

Then  $\lfloor \mathcal{U}_k(\pi) \rfloor$  is an upper bound on the number of breakpoints of  $\pi$  that can be eliminated in  $k$  reversals.

PROOF. The proof is similar to that of Theorem 2.

Given an allotment of  $k$  reversals, we want to eliminate as many breakpoints as possible. Suppose we consider all  $k'$ -moves for  $k' \leq k + 1$ . Packing these  $k'$ -moves into the  $k$  reversals corresponds to finding a matching  $M$  of  $G^{(k+1)}$  satisfying

$$\sum_{e \in M} (|e| - 1) \leq k,$$

since, by Lemmas 4 and 6, each edge  $e$  of  $M$  represents  $|e| - 1$  reversals. The number of breakpoints eliminated by these reversals is

$$\sum_{e \in M} |e|.$$

Having used  $\sum_{e \in M} (|e| - 1)$  reversals, we have  $k - \sum_{e \in M} (|e| - 1)$  remaining in our allotment. Notice that any unconsidered  $k'$ -move, which must have  $k' \geq k + 2$ , will not completely fit in our allotment of  $k$  reversals, since by Lemma 4 such a move takes at least  $k + 1$  reversals. Lemma 4 also implies that any prefix of such a  $k'$ -move that is packed into our allotment, can on average remove at most one breakpoint per reversal. Thus, the number of breakpoints eliminated by the reversals remaining in our allotment, is at most

$$k - \sum_{e \in M} (|e| - 1).$$

Adding the number of breakpoints eliminated by the reversals in  $M$  to the number of breakpoints eliminated by the remaining reversals, the total number eliminated within  $k$  reversals is at most

$$k + \sum_{e \in M} 1.$$

Expressing this matching problem as an integer program, and extending its domain to the reals, results in the linear programming problem of the theorem.  $\square$

For fixed  $k$ , the time and space to compute upper bound  $\mathcal{U}_k$  is the same as for lower bound  $\mathcal{L}_k$ , which is  $O(L(n, n))$  time and  $O(n)$  space.

Given that we can find a series of  $k$  reversals that removes the most breakpoints, how should we piece together a solution from such a series? One extreme is to perform only the first reversal of the series, arrive at a new permutation, and again look ahead  $k$  reversals. The other extreme is to execute all  $k$  of the series. We call the number of reversals that are performed from a series, the *follow-through* of the algorithm.

It might be predicted that a follow-through of one reversal is best, since this retains the maximum flexibility. In our experience, however, this performed unexpectedly poorly, in the sense that the final solution tended to degrade as the look-ahead was increased beyond some critical value. An experimental analysis of this phenomenon would be interesting, but is beyond the scope of our paper.

We also remark that looking ahead farther does not in general guarantee a better solution. After looking ahead  $k + 1$  reversals, we are simply in some state, and unless the permutation can be sorted in  $k + 2$  reversals, this state does not necessarily lead to a shorter solution than the one we arrive at after looking ahead  $k$  reversals.

Nevertheless, choosing a follow-through equal to the look-ahead  $k$  did have the desired property that the quality of the solution tended to improve as the look-ahead was increased, and this is the value that we chose for the experiments of the next section. A follow-through of  $k$  reversals also has the advantage of reducing the number of invocations of the branch-and-bound procedure. It would be interesting, however, to investigate other values, such as a follow-through of  $k/2$  with a look-ahead of  $k$ .

To summarize, our exact algorithm constructs a solution conceptually in three stages, the first two of which are interleaved. The first stage runs the greedy algorithm to lower bound the number of breakpoints that can be eliminated within the look-ahead. The second stage improves this greedy solution by branch-and-bound to a fixed depth, to obtain a series that is optimal within the look-ahead. Successive locally optimal series are then concatenated, to obtain a solution that upper bounds the global problem. The third and final stage improves this solution to a global optimum by a full branch-and-bound computation, now that a good upper bound is in hand.

This bootstrapping approach has proven to be quite effective, as is discussed in the next section.

**4. Computational Results.** To examine the effectiveness of these ideas, we tested a full implementation of the exact and approximation algorithms on biological and simulated data. The implementation comprises approximately 9500 lines of C, of which roughly 2500 lines are a sparse linear programming package.

An unusual aspect of the code is the manipulation of the bidirected graphs of Section 3.2. These graphs are sufficiently different from standard directed and undirected graphs to make the correct implementation of simple computations like shortest paths and cycle enumeration surprisingly tricky. Without going into detail, we note that a straightforward translation of the standard breadth-first search algorithm for single-source shortest paths is not correct, since a vertex can be reached in two different ways from the source: once by paths that end with an in-edge, and once by paths that end in an out-edge. This means that essentially two distances must be maintained for a vertex: an in-edge distance, and an out-edge distance.

We also note that the code enumerates cycles of length at most  $k$  using the method outlined in Section 3.2, which can spend time exponential in  $k$  between

**Table 1.** The number of permutations on  $n$  elements at distance  $d$  from the identity.

$d$	$n$						
	2	3	4	5	6	7	8
0	1	1	1	1	1	1	1
1	1	3	6	10	15	21	28
2		2	15	51	127	263	483
3			2	56	390	1,562	4,635
4				2	185	2,543	16,445
5					2	648	16,615
6						2	2,111
7							2

reporting cycles. The algorithm of Johnson [12] for enumerating the cycles of a directed graph spends time linear in the size of the graph between reporting each cycle, and it would be interesting to see whether this algorithm can be adapted to our application.

In all, we tested the implementation on four types of data: all permutations on a fixed number of elements, published gene order data from the biology literature, random permutations, and permutations generated by scrambling the identity with a fixed number of random reversals.

In the first set of experiments, which served as a good test of the program, we ran the implementation to optimality on all permutations of up to eight elements. The exact distribution of reversal distance from these tests is given in Table 1. Notice that the data supports Golan's conjecture: for each  $n > 2$ , there are exactly two permutations requiring  $n - 1$  reversals (and these are  $\gamma^{(n)}$  and its inverse).

On these experiments, we also measured the worst-case performance ratio of the approximation algorithm. On permutations of up to eight elements, the maximum ratio is  $8/5$ , which is achieved on permutation (4 7 2 6 8 5 3 1).

As an illustration of the algorithm on biological data, we give the example of Figure 10. This permutation gives the order of the 36 genes that are common to the linearized mitochondrial genomes of mammals [3] and the flatworm *Ascaris suum* [25]. While we have been able to solve other permutations arising from mitochondria data in a small amount of time,<sup>16</sup> this 36-element permutation has proven extremely difficult to solve to optimality. The near-optimal solution of the figure was found after 24 s of computation on a 33 Mhz Silicon Graphics Iris

<sup>16</sup> Computing the inversion distance between the mitochondrial genomes of mammals and the yeast *Schizosaccharomyces pombe* took 3.5 min (29 genes and 19 inversions), and between mammals and the fly *Drosophila yakuba* took 1.7 min (37 genes and 16 inversions). We make no claim for the biological significance of the particular solutions found, though the inversion distance tends to reflect evolutionary divergence.

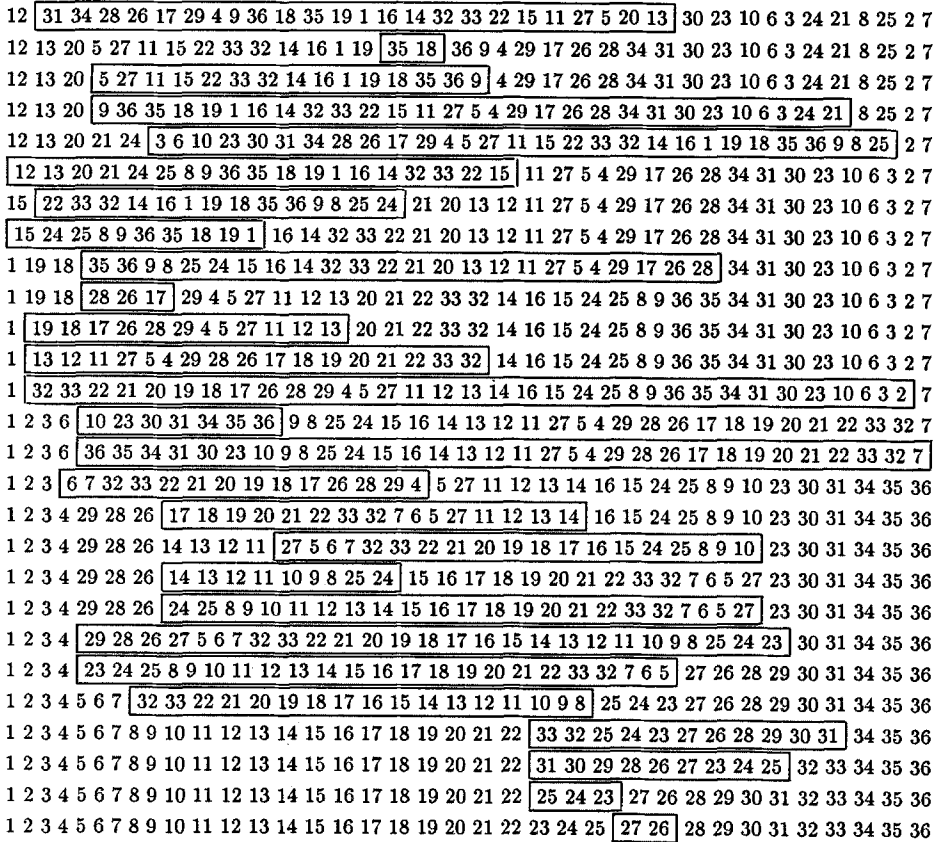


Fig. 10. Near-optimal solution for gene orders from the mitochondrial genomes of mammals and the flatworm *Ascaris suum*. This solution of 27 reversals is provably within two reversals of optimal.

4D/300GTX. The lower and upper bounds from this run, of 25 and 27 reversals, were found using a lower-bound family of six and an upper-bound look-ahead of five reversals.<sup>17</sup> The search tree during look-ahead had a maximum size of 2234 nodes, where this counts all nodes at which a linear programming problem was solved (including nodes that were pruned by the lower bound). From the difference between the upper and lower bounds, we know that this solution is within two reversals of optimal. A family of ten and a look-ahead of eight, which required a search of 408,653 nodes and terminated after 7.5 h of computation, failed to improve the bounds by one reversal.

The limit of what we can reliably solve to optimality is around 30 elements. Table 2 gives the running time and search-tree size to solve a sample of 10 random

<sup>17</sup> As described in Section 3.3, the follow-through for the upper bound was equal to the look-ahead in all experiments.



**Table 2.** Running time and search-tree size for exact solution of random permutations on 30 elements.\*

Running time				Tree size							
				Upper bound algorithm				Exact algorithm			
Min.	Max.	Med.	Ave.	Min.	Max.	Med.	Ave.	Min.	Max.	Med.	Ave.
0:05	43:48	1:29	6:53	228	372,563	11,419	45,393	0	8,291	29	843

\* The lower-bound family is eight, the upper-bound look-ahead is six reversals, and the sample size is ten permutations. The running time is the total time for the exact algorithm in minutes, which includes the time to compute the upper bound. The tree size for the upper-bound algorithm is the maximum size of the search tree during look-ahead. The average reversal distance of the sample is 20.6.

permutations on 30 elements to optimality. Though the lower-bound family and upper-bound look-ahead were the same across these runs, the execution time and search-tree size varied by three orders of magnitude. Since the average running time and tree size are skewed by these outliers, we also give median values, which are more representative of the sample.

That we could solve these instances to optimality is due to the tightness of the bounds. Over the ten permutations, the maximum difference between the upper and lower bound was two reversals, and in each case, the exact solution value was equal to the lower bound. Except for the one permutation on which the exact algorithm examined over 8000 nodes, a common picture emerged from these runs. Optimal solution involved the upper-bound algorithm exploring a rather large tree to find a solution within one or two reversals of optimal; with this solution in hand, the exact algorithm hugged the left chain of its search tree to make up the difference to the lower bound, which was exact.

The effect of varying the family and look-ahead is shown in Table 3. On random permutations of 50 elements, the average lower bound reached a maximum of roughly 36 reversals at family 6, and the average upper bound did not improve much on 38 reversals beyond look-ahead 5.<sup>18</sup> Consequently, these values, family 6 and look-ahead 5, were used in the remaining experiments.

The search tree for these runs was limited to 50,000 nodes. Once this limit was exceeded, the best series known within the look-ahead was used to form the upper bound. Column *T* gives the median tree-size for the sample, where the tree size of a run is the maximum size of the search tree during look-ahead. For a look-ahead of six or more, the majority of runs had a tree size meeting the limit. Column *C* gives the median number of cycle sets on *k* or few breakpoints, which is equivalent to the number of variables in the linear program for the lower bound at family *k*. The number of constraints on these variables is essentially the number of breakpoints in the permutation.

In the third set of experiments, we studied the quality of the approximation algorithm and the upper-bound algorithm on random permutations. Results are

<sup>18</sup> Notice that the average upper bound actually increased at look-ahead 8.

**Table 3.** Lower bound  $L$ , upper bound  $U$ , number of cycle sets  $C$ , and tree size  $T$ , at various families and look-aheads, for random permutations on 50 elements.\*

$k$	$L$		$U$		$C$		$T$	
	Ave.	Dev.	Ave.	Dev.	Med.	Max.	Med.	Max.
1	26.0		43.7		0		917	
2	32.1		40.8		4		2,907	
3	34.5		39.5		14		4,312	
4	35.5		39.1		40		24,242	
5	35.7		38.4		119		36,676	
6	35.9		38.2		353		50,000	
7	35.9		38.0		910		50,000	
8	35.9		38.9		2,464		50,000	

\* The maximum tree-size is 50,000 nodes, and sample size is 10 permutations. Row  $k$  gives the lower bound for family  $k$ , the upper bound for look-ahead  $k$ , the number of cycle sets of at most  $k$  breakpoints for the lower bound, and the maximum tree-size for the upper bound during look-ahead.

given in Table 4, where Dev. denotes the standard deviation. Note that the maximum difference between the upper and lower bound, which is what limits the range of optimal solution, was at most two reversals for  $n$  up to 30, while the average difference between the bounds was around 2.5 reversals for  $n$  up to 50. This suggests that, while we can find optimal solutions for at most around 30 elements, we can find near-optimal solutions that may be acceptable quality for up to 50 elements.

**Table 4.** Lower bound  $L$ , upper bound  $U$ , approximation  $A$ , and tree size  $T$  for random permutations on  $n$  elements.\*

$n$	$L$		$U$		$U - L$		$A$		$A/L$		$T$
	Ave.	Dev.	Ave.	Dev.	Ave.	Max.	Ave.	Dev.	Ave.	Max.	Med.
10	6.0	1.1	6.0	1.1	0.0	0	6.1	1.0	1.02	1.20	0
20	12.6	1.2	13.2	1.5	0.6	2	14.8	1.5	1.18	1.31	517
30	20.8	0.9	21.8	1.6	1.0	2	24.9	1.5	1.20	1.25	4,615
40	28.5	1.1	30.3	1.6	1.8	4	33.8	2.0	1.19	1.24	23,712
50	35.9	1.4	38.4	1.8	2.5	5	43.7	2.7	1.22	1.31	36,676
60	43.6	1.1	46.7	1.5	3.1	5	52.6	1.6	1.21	1.26	50,000
70	51.7	1.3	56.7	2.4	5.0	8	63.4	3.0	1.23	1.30	50,000
80	58.9	1.0	64.5	2.0	5.6	8	72.3	2.3	1.23	1.26	50,000
90	67.6	1.4	74.1	2.1	6.5	8	83.2	1.8	1.23	1.25	50,000
100	74.2	1.1	82.4	2.6	8.2	10	91.9	2.7	1.24	1.27	50,000

\* The lower-bound family is six, the upper-bound look-ahead is five reversals, the maximum tree size is 50,000 nodes, and the sample size is ten permutations. Approximation  $A$  is the number of reversals from the greedy algorithm. Tree size  $T$  is the maximum tree-size for the upper bound during look-ahead.

**Table 5.** Upper bound  $U$  and lower bound  $L$  for permutations on 100 elements generated by  $k$  random reversals.\*

$k$	$U$		$L$	
	Ave.	Dev.	Ave.	Dev.
5	5.0	0.0	5.0	0.0
10	10.0	0.0	10.0	0.0
15	15.2	0.4	15.0	0.0
20	19.5	0.9	19.4	1.0
25	25.0	0.9	24.5	0.9
30	30.0	1.3	29.3	1.0
35	34.6	2.1	33.5	1.8
40	39.0	1.6	37.3	2.0
45	43.9	2.2	42.5	2.2
50	47.3	2.0	45.5	1.4
55	51.8	3.0	48.5	2.8
60	54.2	2.7	50.7	2.0
65	58.5	2.6	54.3	2.5
70	60.0	4.3	56.0	3.1
75	62.0	2.8	57.7	2.4
80	64.9	3.9	60.4	2.8
85	67.0	2.6	61.9	2.0
90	69.8	3.0	64.9	2.2
95	72.2	2.2	65.1	1.5
100	72.0	2.9	65.3	1.7

\* The upper-bound look-ahead is five reversals, the lower-bound family is six, the maximum tree-size is 50,000 nodes, and the sample size is ten permutations.

The average performance ratio of the approximation algorithm for the sample was around  $5/4$ , while the poorest ratio was less than  $4/3$ .

In the final set of experiments, we were interested in the following question: How well does reversal distance recover the actual number of reversals performed on a permutation? To examine this, we generated permutations by scrambling the identity with  $k$  random reversals, taking care not to reverse single elements. Table 5 gives upper and lower bounds on the reversal distance, for permutations on 100 elements with up to 100 random reversals. For  $k < \frac{1}{2}n$ , the average upper bound estimated  $k$  to within one reversal, but for  $k \geq \frac{1}{2}n$ , it increasingly underestimated  $k$ .

Notice that when both endpoints of every reversal in a series of length  $k$  fall at new positions, the reversal distance for the resulting permutation is precisely  $k$ .<sup>19</sup> Since a permutation on  $n$  elements can have at most  $n + 1$  breakpoints, the maximum number of reversals for which this can happen is  $\frac{1}{2}(n + 1)$ , which may partly explain why we observe a change in behavior around  $k = \frac{1}{2}n$ .

<sup>19</sup> This follows because  $k$  is a lower bound, as well as an upper bound, on the reversal distance: the number of breakpoints divided by two, which lower bounds reversal distance, is in this situation equal to  $k$ .

**5. Conclusion.** Algorithmic study of the reversal distance between permutations is in its earliest stages. We have presented two algorithms: a greedy approximation algorithm, and a branch-and-bound exact algorithm. By analyzing the  $\Delta\Phi$ -sequence of the greedy algorithm, we were able to show it achieves a worst-case approximation factor of 2, and by applying matchings, shortest paths, and linear programming, we derived a class of nontrivial lower bounds for our exact algorithm.

Experiments with the exact algorithm indicate that we can solve random permutations up to around 30 elements to optimality, usually in a few minutes, and that for permutations generated by  $k$  random reversals, the average solution value is a good estimate of the number of reversals when  $k$  is less than half the number of elements.

We close with some conjectures and lines for future research.

*5.1. Further Research.* One question that we would like to resolve is the computational complexity of sorting by reversals. We conjecture it is NP-complete, and believe the crux of the problem is the following.

CONJECTURE 1. *Deciding whether  $\text{OPT}(\pi) \leq \frac{1}{2}\Phi(\pi)$  is NP-complete.*

Notice that this special case has a lot of additional structure—for instance, every breakpoint must be paired with another in  $G(\pi)$ —which should simplify a proof (or disproof). On the other hand, a disproof by an efficient algorithm would improve our lower bounds of Section 3.2.

To design a more efficient exact algorithm, we need theorems on the structure of a solution. In particular, do we really need to consider all  $\binom{n}{2}$  reversals to find an optimal solution?

It is natural to think we can throw out reversals that cut strips, since such a reversal separates elements that will have to be joined together later. Unfortunately, this is incorrect. Permutation (3 4 1 2), for example, requires three reversals if we do not cut strips, yet it can be sorted in two reversals, as follows:

$$|3\ 4|1\ 2| \vdash 1|4\ 3\ 2| \vdash 1\ 2\ 3\ 4.$$

Nevertheless, we believe the following is true.

CONJECTURE 2. *Every permutation has an optimal solution that does not cut strips other than at their first or last element.*

If true, this reduces the number of candidate reversals from  $O(n^2)$  to  $O(\Phi^2(\pi))$ . Our belief in Conjecture 2 is based on the following.

CONJECTURE 3. *Every permutation has an optimal solution that never increases the number of breakpoints.*

Another tempting idea is to decompose the permutation into knots, where a *knot* of  $\pi$  is an interval  $[a, b]$  such that

- $\pi(i) \in [1, a)$  for  $i \in [1, a)$ ,
- $\pi(i) \in [a, b]$  for  $i \in [a, b]$ , and
- $\pi(i) \in (b, n]$  for  $i \in (b, n]$ ,

and then solve knots independently. This is appealing, since a permutation can always be sorted without an element crossing a knot. Perhaps surprisingly, it is not optimal. Permutation (2 1 5 6 3 4), for example, has knots  $[1, 2]$  and  $[3, 6]$ , and takes four reversals if we solve the knots separately. Yet, by working on the knots together, the permutation can be sorted in three reversals, as shown below:

$$|2\ 1|5\ 6|3\ 4| \vdash |6\ 5|1\ 2\ 3\ 4| \vdash |4\ 3\ 2\ 1|5\ 6 \vdash 1\ 2\ 3\ 4\ 5\ 6.$$

Finally, we note that the suboptimality of the greedy algorithm is not due to the fact that it chooses arbitrarily among moves that appear equally good. It is necessary to consider 1-moves when 2-moves are available. Permutation (5 6 2 1 8 7 3 4), for example, has one 2-move,  $[3, 6]$ , and taking this move leads to a solution of four reversals. Yet the permutation can be sorted in three reversals, as follows:

$$|5\ 6|2\ 1|8\ 7|3\ 4| \vdash 1\ 2|6\ 5|8\ 7|3\ 4| \vdash 1\ 2|6\ 5\ 4\ 3|7\ 8 \vdash 1\ 2\ 3\ 4\ 5\ 6\ 7\ 8.$$

Before concluding, we remark that the biological motivation for reversal distance suggests several variations. Here we identify five aspects of the problem that may be varied:

- (1) The *linear/circular variation*. Some organisms and organelles have circular chromosomes, for which the data is the order of genes around a circle. In this case, cyclic shifts of a permutation are isomorphic, which means  $(\pi_1 \cdots \pi_n)$  is equivalent to  $(\pi_i \cdots \pi_n \pi_1 \cdots \pi_{i-1})$ . A solution is then a series of reversals and cyclic shifts that transform one permutation to another, and we seek a solution with the fewest reversals.
- (2) The *signed/unsigned variation*. A gene is not a point on a line or circle; it is a region of sequence. The sequence has a reading direction, and a reversal reverses not only the order of genes, but their direction as well.

When two chromosomes are compared, information may be available on gene direction, as well as order. We can encode this information by associating a sign with each element: positive for the forward direction, and negative for the reverse. A reversal then changes the sign of elements it reverses.

In the signed case, we have two signed permutations, and we seek a shortest series of signed reversals that transform one into the other. In some respects, this case is easier to analyze. Conjecture 3, for instance, holds.

- (3) The *directed/undirected variation*. When a chromosome is analyzed, sometimes only the adjacency of genes is determined, not their absolute order. In the

linear case, this means the given order may be left to right, or right to left, and in the circular case, clockwise or counterclockwise.

- (4) The *weighted/unweighted variation*. We have assumed that all reversals are equally probable, so it is appropriate to count their number. For a first approximation this is reasonable, but a finer analysis might weight reversals by a function of their length, and find a series of minimum total weight.
- (5) The *pairwise/multiple variation*. Even though the objective of the reversal distance problem is not to find an alignment of the input, as is often the case with sequence comparison, we can define a multiple-chromosome comparison problem in analogy to multiple alignment. Given a set of orderings of the same genes, how can we infer an evolutionary tree, and a set of hypothetical ancestral orderings, so as to minimize the total reversal distance of the tree? Admittedly, a practical exact solution method appears unlikely, but a good approximation method would be useful as well.

Clearly, there are many possibilities for exploration.

**Acknowledgments.** Guillaume Leduc, in analyzing signed reversals, contributed to the line of proof that  $\Phi(\pi)$  bounds  $\text{OPT}(\pi)$ . Holger Golan began the study of the distribution of  $\text{OPT}(\pi)$  for small  $n$ , and with his conjectures on the diameter of the symmetric group generated by reversals, stimulated the search for better lower bounds. Michel Berkelaar generously provided the sparse linear programming code used in our experiments.

Our thanks to the referees for their comments. The first author also wishes to thank Dan Gusfield, Paul Stelling, and Lucas Hui for many helpful discussions.

*Note Added in Proof.* A shorter version of this paper appeared as [12a], and an extension to signed permutations is given in [12b]. For signed permutations, we observe an average difference, between the greedy approximation and a simplified lower bound, of less than one reversal, for  $n$  up to 10,000. Bafna and Pevzner [2a] have since improved the performance ratio to  $7/4$  for unsigned permutations, and  $3/2$  for signed permutations, and have also established Gollan's conjecture.

## References

- [1] Aigner, M., and D. B. West. Sorting by insertion of leading elements. *Journal of Combinatorial Theory, Series A*, **45**, 306–309, 1987.
- [2] Amato, N., M. Blum, S. Irani, and R. Rubinfeld. Reversing trains: a turn of the century sorting problem. *Journal of Algorithms*, **10**, 413–428, 1989.
- [2a] Bafna, V., and P. A. Pevzner. Genome rearrangements and sorting by reversals. *Proceedings of the 34th Symposium on Foundations of Computer Science*, November 1993, pp. 148–157.
- [3] Bibb, M. J., R. A. van Etten, C. T. Wright, M. W. Walberg, and D. A. Clayton. Sequence and gene organization of mouse mitochondrial DNA. *Cell*, **26**, 167–180, 1981.

- [4] Dobzhansky, T. *Genetics of the Evolutionary Process*. Columbia University Press, New York, 1970.
- [5] Driscoll, J. R., and M. L. Furst. Computing short generator sequences. *Information and Computation*, **72**, 117–132, 1987.
- [6] Even, S., and O. Goldreich. The minimum-length generator sequence problem is NP-hard. *Journal of Algorithms*, **2**, 311–313, 1981.
- [7] Furst, M., J. Hopcroft, and E. Luks. Polynomial-time algorithms for permutation groups. *Proceedings of the 21st Symposium on Foundations of Computer Science*, 1980, pp. 36–41.
- [8] Garey, M. R., and D. S. Johnson. *Computers and Intractability: A Guide to The Theory of NP-Completeness*. Freeman, New York, 1979.
- [9] Gates, W. H., and C. H. Papadimitriou. Bounds for sorting by prefix reversal. *Discrete Mathematics*, **27**, 47–57, 1979.
- [10] Golan, H. Personal communication, 1991.
- [11] Jerrum, M. R. The complexity of finding minimum-length generator sequences. *Theoretical Computer Science*, **36**, 265–289, 1985.
- [12] Johnson, D. B. Finding all the elementary circuits of a directed graph. *SIAM Journal on Computing*, **4**(1), 77–84, 1975.
- [12a] Kececioglu, J., and D. Sankoff. Exact and approximation algorithms for the inversion distance between two chromosomes. *Proceedings of the 4th Symposium on Combinatorial Pattern Matching*, Lecture Notes in Computer Science, vol. 684, Springer-Verlag, Berlin, June 1993, pp. 87–105. (An earlier version appeared as “Exact and approximation algorithms for sorting by reversals,” Technical Report 1824, Centre de recherches mathématiques, Université de Montréal, July 1992).
- [12b] Kececioglu, J., and D. Sankoff. Efficient bounds for oriented chromosome-inversion distance. *Proceedings of the 5th Symposium on Combinatorial Pattern Matching*, Lecture Notes in Computer Science, vol. 807, Springer-Verlag, Berlin, June 1994, pp. 307–325.
- [13] Knuth, D. E. *The Art of Computer Programming*, Vol. 3. Addison-Wesley, Reading, MA, 1973.
- [14] Mannila, H. Measures of presortedness and optimal sorting algorithms, *IEEE Transactions on Computers*, **34**, 318–325, 1985.
- [15] Micali, S. and V. Vazirani. An  $O(\sqrt{|V|} \cdot |E|)$  algorithm for finding maximum matchings in general graphs. *Proceedings of the 21st Symposium on Foundations of Computer Science*, 1980, pp. 17–27.
- [16] Nadeau, J. H., and B. A. Taylor. Lengths of chromosomal segments conserved since divergence of man and mouse. *Proceedings of the National Academy of Sciences of the USA*, **81**, 814, 1984.
- [17] O’Brien, S. J., ed. *Genetic Maps: Locus Maps of Complex Genomes*. 6th edition. Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, 1993.
- [18] Palmer, J. D., B. Osorio, and W. F. Thompson. Evolutionary significance of inversions in legume chloroplast DNAs. *Current Genetics*, **14**, 65–74, 1988.
- [19] Sankoff, D., G. Leduc, N. Antoine, B. Paquin, B. F. Lang, and R. Cedergren. Gene order comparisons for phylogenetic inference: evolution of the mitochondrial genome. *Proceedings of the National Academy of Sciences of the USA*, **89**, 6575–6579, 1992.
- [20] Schöniger, M., and M. S. Waterman. A local algorithm for DNA sequence alignment with inversions. *Bulletin of Mathematical Biology*, **54**, 521–536, 1992.
- [21] Sessions, S. K. Chromosomes: molecular cytogenetics. In *Molecular Systematics*, D. M. Hillis and C. Moritz, eds., Sinauer, Sunderland, MA, 1990, pp. 156–204.
- [22] Tichy, W. F. The string-to-string correction problem with block moves. *ACM Transactions on Computer Systems*, **2**(4), 309–321, 1984.
- [23] Wagner, R. A. On the complexity of the extended string-to-string correction problem. In *Time Warps, String Edits, and Macromolecules: The Theory and Practice of Sequence Comparison*, D. Sankoff and J. B. Kruskal, eds., Addison-Wesley, Reading, MA, 1983, pp. 215–235.
- [24] Watterson, G. A., W. J. Ewens, T. E. Hall, and A. Morgan. The chromosome inversion problem. *Journal of Theoretical Biology*, **99**, 1–7, 1982.
- [25] Wolstenholme, D. R., J. L. MacFarlane, R. Okimoto, D. O. Clary, and J. A. Wahleithner. Bizarre tRNAs inferred from DNA sequences of mitochondrial genomes of nematode worms. *Proceedings of the National Academy of Sciences of the USA*, **84**, 1324–1328, 1987.