

Sorting Signed Permutations by Reversals in Nearly-Linear Time*

Bartłomiej Dudek¹, Paweł Gawrychowski¹, and Tatiana Starikovskaya²

¹Institute of Computer Science, University of Wrocław, Poland
¹{bartlomiej.dudek,gawry}@cs.uni.wroc.pl

²DIENS, École normale supérieure, PSL Research University, France
¹tat.starikovskaya@gmail.com

Abstract

Given a signed permutation on n elements, we need to sort it with the fewest reversals. This is a fundamental algorithmic problem motivated by applications in comparative genomics, as it allows to accurately model rearrangements in small genomes. The first polynomial-time algorithm was given in the foundational work of Hannenhalli and Pevzner [J. ACM'99]. Their approach was later streamlined and simplified by Kaplan, Shamir, and Tarjan [SIAM J. Comput.'99] and their framework has eventually led to an algorithm that works in $\mathcal{O}(n^{3/2}\sqrt{\log n})$ time given by Tannier, Bergeron, and Sagot [Discr. Appl. Math.'07]. However, the challenge of finding a nearly-linear time algorithm remained unresolved. In this paper, we show how to leverage the results on dynamic graph connectivity to obtain a surprisingly simple $\mathcal{O}(n \log^2 n / \log \log n)$ time algorithm for this problem.

1 Introduction

The main goal of comparative genomics is to design efficient methods for comparing genomes of different species. A genome can be thought of as a linear sequence of genes, identified by numbers, with every gene having a direction, or a sign. Already in the eighties, it has been observed that seemingly very different genomes might actually have a small reversal distance, defined as the smallest number of reversals to transform one genome into another. We refer to the comprehensive treatment of different aspects of such an approach in the book by Pevzner [15]. This brings the question of computing this number efficiently.

A clean combinatorial formulation of this problem was first formulated by Day and Sankoff [7]. We think that each genome is a permutation π on $\{1, 2, \dots, n\}$ where every element additionally has a sign $+$ or $-$. A reversal takes a segment of such a π , reverses it and additionally flips the sign of its elements. Given a signed permutation π and π' , we want to find a shortest sequence of reversals that transforms π into π' . The length of such a sequence is called the reversal distance between π and π' . By renaming the elements, it is enough to consider sorting a given signed permutation π , that is, transforming it into $(1, 2, \dots, n)$.

The challenge of designing an efficient algorithm for sorting signed permutations by reversals has a long history. Initially, only approximation algorithms were known [2, 13]. In their famous

*This work was partially funded by the grant ANR-20-CE48-0001 from the French National Research Agency (ANR).

JACM paper titled “Transforming Cabbage into Turnip”, Hannenhalli and Pevzner [9] were able to design a polynomial time algorithm for this problem. This is somewhat surprising, as the unsigned version of the problem is in fact NP-hard [6]. Their main contribution was a duality theorem that connects the reversal distance with some combinatorial parameters, which results in a fairly complicated algorithm running in $\mathcal{O}(n^4)$ time. Very soon afterwards, Berman and Hannenhalli [9] showed how to improve its running time to $\mathcal{O}(n^2\alpha(n))$, however the underlying combinatorial theory was still very complicated. Fortunately, Kaplan, Shamir, and Tarjan [11] were able to considerably simplify the combinatorial theory and the resulting algorithm, thus obtaining a faster quadratic time algorithm, and Bergeron [4] simplified it even further, obtaining a different algorithm with cubic running time. However, no worst-case subquadratic time algorithm was known at this point. Ozery-Flato and Shamir [14] designed a family of permutations on which any algorithm based on a similar approach needs at least quadratic time, and explicitly raised the question of designing a subquadratic algorithm.

An empirical answer to this question was provided by Kaplan and Verbin [12], who designed an algorithm based on a random walk in a certain graph. On a random input, their algorithm seems to find the correct shortest sequence with high probability in $\mathcal{O}(n\sqrt{n\log n})$ time. More importantly, at the heart of their algorithm lies a structure for maintaining a permutation under reversals in $\mathcal{O}(\sqrt{n\log n})$ time per operation. An extension of this structure was used by Tannier, Bergeron, and Sagot [19] to design a worst-case $\mathcal{O}(n\sqrt{n\log n})$ time algorithm that works for any input. Their algorithm departs in an interesting way from the previous approaches, which were based on repeatedly applying the so-called safe reversals (reversals that decrease the distance by 1). Instead, they maintain a sequence of reversals, and repeatedly increase its length by inserting (possibly multiple) reversals somewhere in the middle (and not just at the end). Their algorithm crucially depend on the observation that the position in which we insert the new reversals never moves to the right in the current sequence, which allows for an efficient implementation. Finally, Han [8] showed how to tweak the structure of Kaplan and Verbin [12] to implement some of its operations in $\mathcal{O}(\sqrt{n})$ time. While not fully described in his paper, it is plausible that such an approach in fact works for speeding up all the operations.

After the subquadratic time algorithm of Tannier, Bergeron, and Sagot [19], the next challenge is understanding if the complexity of sorting signed permutations by reversals is $\tilde{\mathcal{O}}(n^{1.5})$, or perhaps is there a near-linear time algorithm? On a more applied side, Swenson, Rajan, Lin, and Moret [18] identified a data-dependent parameter k , observed that it appears to be a small constant for a random input, and designed an $\mathcal{O}(n\log n + kn)$ time algorithm. However, no worst-case near-linear time algorithm was known. Interestingly, computing the reversal distance itself can be done in linear time [1, 5]. However, constructing the corresponding sequence seems more challenging.

Our result. We design a near-linear time for sorting a signed permutation with reversals. More specifically, our algorithm outputs a shortest sequence of reversals sorting a given signed permutation on n elements in near-linear time. First, we show how to combine the algorithm of Tannier, Bergeron, and Sagot [19] with a structure for maintaining a graph under inserting/deleting edges and maintaining a minimum spanning tree. This immediately gives us an $\mathcal{O}(n\log^4 n)$ time algorithm by plugging in e.g. the structure of Holm et al. [10]. Next, we show that in fact maintaining any spanning forest is sufficient, which allows for a faster $\mathcal{O}(n\log^2 n)$ time algorithm [10], or $\mathcal{O}(n\log^2 n/\log\log n)$ by plugging in the fastest deterministic structure [20] and some careful bookkeeping.

Our method. Our algorithm builds on the work of Tannier, Bergeron, and Sagot [19]. The high-level description of their algorithm is as follows. They maintain a sequence S of reversals, partitioned into S_1 and S_2 . S_2 will be the suffix of the final sequence, but there might be some further reversals inserted into S_1 . They maintain the current signed permutation obtained by applying all the reversals in S_1 to the initial permutation. They also maintain the set of remaining reversals V that can create a new adjacency, that is, bring elements i and $i + 1$ to the adjacent positions. Then, as long as there exists a reversal in V such that its endpoints contain elements with different signs in the current permutation, such a reversal is applied, removed from V , and appended to S_1 . Otherwise, the last reversal of S_1 is undone and prepended to S_2 . The main difficulty is to detect a reversal in V such that its endpoints contain elements with different signs, given that we keep reversing different segments of the permutation (recall that this involved flipping the sign of every element in the segment). Instead of designing a new structure for this problem, we reduce it to the well-known fully dynamic connectivity problem. More specifically, our graph consists of two paths on nodes $\{0, 1, \dots, n + 1\}$ and $\{0', 1', \dots, (n + 1)'\}$. We think that the paths consist of blue edges. Throughout the execution of the algorithm, we maintain the two paths, and for each $i = 0, 1, \dots, (n + 1)'$ the node i is on one path while the node i' is on the other path. The order on the nodes on each path corresponds to the current permutation (disregarding the signs), and the path starting at node 0 contains the node i if and only if the element i of the permutation has been flipped an even number of times. It is easy to see how to maintain such invariants by changing 4 edges per reversal in the permutation. Next, we also encode every reversal in V as two red edges in the graph. This is done in such a way that the endpoints of a reversal in V contain elements with different signs if and only if its corresponding red edges connect different blue paths in the current graph. This does not require any modifications to the graph after a reversal, except that we need to remove its corresponding red edges when removing a reversal from V . We maintain a fully dynamic connectivity structure for the graph; we stress that, even though we have edges of two colours, we simply maintain all of them in the structure, disregarding the colours. Now, checking if there exists a reversal in V such that its endpoints contain elements with different signs reduces to checking if the whole graph is connected. Extracting such a reversal requires a bit more work, as we have edges of two colors in our graph, and we need to find a red edge in the current spanning tree. We show how to implement this efficiently by binary searching over the unique path connecting nodes 0 and $0'$ in the spanning tree.

2 Preliminaries

In this section, we describe the previous framework and summarise it as a concise interface on signed permutations in Theorem 2.7. Our result follows from an efficient implementation of this interface presented in Section 3.

We are mostly using the naming convention and definitions from [19]. A signed permutation π on n elements is a sequence where each element of $\{1, 2, \dots, n\}$ appears exactly once and has a sign, $+$ (often omitted) or $-$. By π_i , we denote the i th element of the sequence and by π_i^{-1} an index j such that $\pi_j = \pm i$. In the problem of sorting by reversals, we are given a signed permutation π and must find a shortest sequence of reversals transforming π into the identity permutation $Id = (1, 2, \dots, |\pi|)$. A reversal ρ of an interval $[i, j]$ transforms a signed permutation $\pi = (\pi_1, \pi_2, \dots, \pi_n)$ into a signed permutation $\pi \cdot \rho = (\pi_1, \dots, \pi_{i-1}, -\pi_j, -\pi_{j-1}, \dots, -\pi_i, \pi_{j+1}, \dots, \pi_n)$. The length of a shortest sequence of reversals sorting π is denoted by $d(\pi)$.

Consider a signed permutation π on n elements. Following the convention, we extend the permutation π by $+0$ at the beginning and $+(n + 1)$ at the end, that is hereafter we consider only

permutations $\pi = (+0, \pi_1, \dots, \pi_n, +(n+1))$. We will never involve 0 or $n+1$ in a reversal and therefore their signs will never change.

Recall the definition of the overlap graph of a signed permutation π , $OV(\pi)$, introduced by Kaplan, Shamir, and Tarjan [11]. See Figure 1. We associate two points π_i^- and π_i^+ to every $i \in [n]$, and also 0^+ to 0 and $(n+1)^-$ to $(n+1)$. These points are linearly ordered in such a way that $\pi_i^- \prec \pi_i^+$ if $\pi_i \geq 0$ and $\pi_i^+ \prec \pi_i^-$ otherwise, and $\pi_i^x \prec \pi_j^y$ whenever $i < j$ for any combination of $x, y \in \{-, +\}$. We then add $n+1$ arcs in such a way that every point is an endpoint of one arc, namely, an arc v_i connects i^+ and $(i+1)^-$. The *overlap graph* $OV(\pi)$ of a permutation π is a graph in which the nodes correspond to the $n+1$ arcs v_i , and the nodes corresponding to two arcs are connected by an edge if the intervals that the arcs span intersect, but none is contained in the other. A node v_i is *black* if i and $(i+1)$ have different signs in π and *white* otherwise.

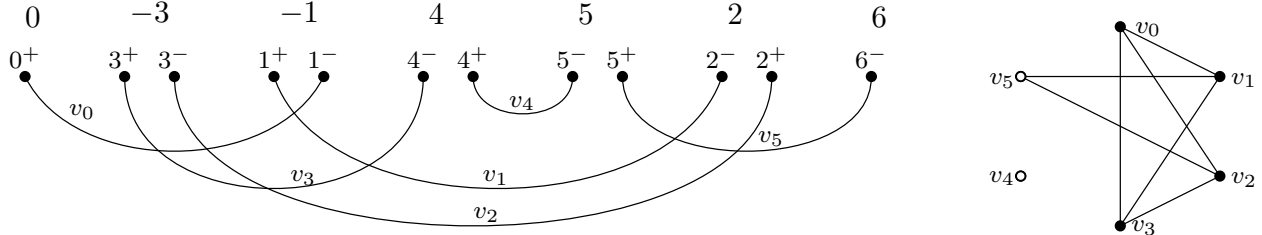


Figure 1: Left: permutation $\pi = (0, -3, -1, 4, 5, 2, 6)$ with its associated points π_i^+ and π_i^- and arcs v_i . Right: the overlap graph $OV(\pi)$.

Fact 2.1 ([19]). *The isolated nodes of $OV(\pi)$ are white and correspond to elements on positions j such that $\pi_j + 1 = \pi_{j+1}$. Such a pair of positions $j, j+1$ is called an adjacency.*

A reversal on a permutation corresponds to a natural transformation of its overlap graph, later referred to as *toggle*. Toggles are defined for general graphs where each node is either black or white, not necessarily overlap graphs. Fix such a graph G . Let $\Gamma(v)$ be the neighborhood of a node v of G and $\Gamma^+(v) = \Gamma(v) \cup \{v\}$. Define $\text{toggle}(v)$ as a local complementation of $\Gamma^+(v)$, that is negating the color of every node and complementing the edges in the subgraph of G induced by $\Gamma^+(v)$ (adding an edge where there was no edge and removing all existing edges). The operation is only allowed for black nodes v . Note that after this operation, v always becomes white and isolated. By G/v we denote the graph obtained from G after applying $\text{toggle}(v)$. We naturally extend this notation to a sequence $S = (s_1, \dots, s_k)$ of nodes of G , namely $G/S := G/s_1/\dots/s_k$.

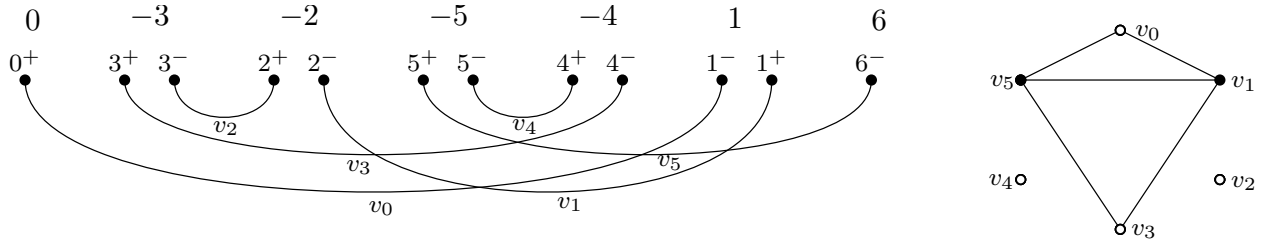


Figure 2: Left: permutation $\pi = (0, -3, -1, 4, 5, 2, 6)$ from Figure 1 after reversal $\rho(v_2)$. Right: the overlap graph $OV(\pi)/v_2 = OV(\pi \cdot \rho(v_2))$, see Lemma 2.2.

Let v be a black node of $OV(\pi)$. Denote by $\rho(v)$ a reversal of an interval consisting of the occurrences of all elements k such that both k^+ and k^- are spanned by the arc v , including its endpoints. Formally, if $v = v_i$, let $a = \min\{\pi_i^{-1}, \pi_{i+1}^{-1}\}$ and $b = \max\{\pi_i^{-1}, \pi_{i+1}^{-1}\}$. If $\pi_a + \pi_b = 1$, then $\rho(v)$ is a reversal of $[a, b - 1]$, and of $[a + 1, b]$ otherwise. See Figure 2.

Lemma 2.2 ([9, 11, 19]). *For a permutation π and a black node v of $OV(\pi)$, we have that $OV(\pi)/v = OV(\pi \cdot \rho(v))$.*

Lemma 2.2 implies that it suffices to find a shortest sequence of toggles that makes all nodes of $OV(\pi)$ white and isolated (extending π by 0 and $n + 1$ at the ends guarantees that the resulting graph is $OV(Id)$).

A difficulty on the way to computing such a sequence is that $OV(\pi)$ might have a non-singleton connected component containing only white nodes (“all-white connected component”) that we will not be able to change by toggling. Fortunately, there is a way to overcome this hurdle:

Theorem 2.3 ([1, 5, 11]). *Given a signed permutation π , it is possible to find in $\mathcal{O}(n)$ time a sequence of t reversals such that when we apply them to π we obtain a permutation π' for which $OV(\pi')$ does not have non-singleton all-white connected components and $d(\pi') = d(\pi) - t$.*

To compute a shortest sequence of toggles efficiently, we will need one more definition:

Definition 2.4 ([19]). *A black node v of a graph G is called safe if there is no non-singleton all-white connected component in G/v .*

Hannenhalli and Pevzner [9] showed that instead of looking for a shortest sequence of toggles making all nodes in $OV(\pi)$ white and isolated, it suffices to find any sequence of toggles of safe nodes:

Theorem 2.5 ([9, 19]). *If v is a safe node of $OV(\pi)$, then $d(\pi \cdot \rho(v)) = d(\pi) - 1$.*

Tannier et al. [19] showed an algorithm that receives a graph G with no non-singleton all-white connected components, and constructs a sequence $S = (v_1, v_2, \dots, v_k)$ such that each node v_{i+1} is safe in $G/v_1/\dots/v_i$, and every node of G/S is white and isolated. The algorithm proceeds in iterations while maintaining a sequence of toggles partitioned into S_1 and S_2 together with G/S_1 , terminating when there are no black nodes in G/S_1 . Every iteration consists of two phases. In the first phase, the algorithm repeatedly toggles black nodes in G/S_1 , appending each toggled node to S_1 . Then, if S_2 is nonempty and its first element is not a black node in G/S_1 , the last element of S_1 is removed and its corresponding toggle undone. In the second phase, as long as there are no black nodes in G/S_1 the last element of S_1 is removed from S_1 and prepended to S_2 and its corresponding toggle undone. See Algorithm 1. There, $()$ is the empty sequence, S_1, S_2 denotes the concatenation of two sequences and S, v or v, S should be read as $S, (v)$ or $(v), S$ respectively. $S[1]$ is the first element of S and $S[2\dots]$ is the sequence containing all elements of S but the first one. Correctness of the algorithm is guaranteed by the following theorem:

Theorem 2.6 ([19]). *Consider a graph G with n black or white nodes with no non-singleton all-white connected components. A sequence v_1, \dots, v_k of nodes in G such that v_i is a black node in $G/v_1/\dots/v_{i-1}$ for $1 \leq i \leq k$ and $G/v_1/\dots/v_k$ has only white isolated nodes can be found by an algorithm that maintains a sequence S of nodes, the graph G/S after applying the sequence of toggles, a subset V of its nodes and performs the following operations:*

1. find a black node in V ,

2. $\text{toggle}(v)$ for a given black node v ,
3. undo the last operation $\text{toggle}(v)$, where v is the last element from S and remove v from the end of S ,
4. remove from V all nodes that became isolated and white after the last toggle operation.

Initially, the set V consists of all non-isolated nodes from G . Every node is toggled at most once (and possibly this is once undone) and each of the above operations is performed in total $\mathcal{O}(n)$ times. Apart from supporting the above operations, the algorithm takes $\mathcal{O}(n)$ time.

Unfortunately, [19] does not provide a complete proof of the above theorem. More specifically, it first shows in Theorem 3 that any maximal but not total valid sequence S for G can be split into $S = S_1, S_2$ so that we can find a nonempty sequence S' for which S_1, S', S_2 is a valid sequence for G . Next, it states Algorithm 1, but without any proof of correctness. At least two issues need to be clarified. First, the algorithm needs to maintain a maximal sequence, so the proof of Theorem 3 needs to be modified so that given a maximal sequence it outputs a longer maximal sequence. Second, the efficiency of the algorithm hinges on the observation that the position at which we split S into S_1 and S_2 never moves to the right, so this needs to be proven. As the algorithm serves as a foundation of our improvement, we felt obliged to provide a full proof, see Appendix A. We stress that the algorithm itself is exactly the same, and the main ideas of the proof were already present in the original paper.

Algorithm 1 Algorithm of Tannier, Bergeron, and Sagot [19]

```

1: function PROCESS(graph  $G$  with no non-singleton all-white components)
2:    $V$  is the set of all nodes of  $G$  that are non-isolated or black
3:    $S_1, S_2 := ()$ 

4:   function UNDOANDREMOVELASTMOVEFROM $S_1()$ 
5:      $w :=$  last element of  $S_1$ 
6:     remove the last element (that is:  $w$ ) from  $S_1$ 
7:     undo  $\text{toggle}(w)$ 
8:     return  $w$ 

9:   while there is a black node  $v$  in  $V$  do
10:     $\text{toggle}(v)$ 
11:    remove from  $V$  all nodes that became isolated and white (in particular: node  $v$ )
12:     $S_1 := S_1, v$ 
13:    if  $S_2 \neq ()$  and  $S_2[1]$  is white then
14:      UNDOANDREMOVELASTMOVEFROM $S_1()$ 
15:    if  $V$  is empty then
16:      return  $S_1, S_2$ 
17:    while there is no black node in  $V$  do
18:       $w :=$  UNDOANDREMOVELASTMOVEFROM $S_1()$ 
19:       $S_2 := w, S_2$ 
20:    go to line 9

```

Now we utilize all the above observations and show that we can efficiently sort by reversals

using only a small set of operations on the permutation. This resembles but extends the interface provided by [8].

Theorem 2.7. *For every signed permutation π of n elements there exists an algorithm that finds a shortest sequence of reversals sorting π . The algorithm maintains π and a set $V \subseteq \{1, 2, \dots, n\}$ under the following operations:*

1. query for π_i or π_i^{-1} ,
2. find $i \in V$ such that i and $i + 1$ have different signs in π ,
3. apply to π a signed reversal of a given interval,
4. remove an element from V .

The algorithm performs $\mathcal{O}(n)$ such operations and additionally takes $\mathcal{O}(n)$ time.

Proof. We start with extending the permutation π with 0 and $n + 1$ at the ends. If $OV(\pi)$ contains non-singleton all-white connected components, we first compute in $\mathcal{O}(n)$ time the sequence of reversals from Theorem 2.3 and apply them to π . From now on, we assume that every all-white connected components of $OV(\pi)$ is a singleton.

Next, we simulate the algorithm from Theorem 2.6 on $OV(\pi)$ based on Lemma 2.2. By Fact 2.1, we initialize $V = \{\max\{\pi_j, -\pi_{j+1}\} \text{ for } j : \pi_j + 1 \neq \pi_{j+1}\}$. By definition, a black node in $OV(\pi)$ corresponds to i such that i and $i + 1$ have different signs in π , so we can find it using operation 2 on π . By Lemma 2.2, toggling a black node v from $OV(\pi)$ corresponds to a reversal of the interval $\rho(v)$. To undo the last toggle $\text{toggle}(v)$, it suffices to reverse $\rho(v)$ one more time as reversing an interval twice is the identity. However, we cannot retrieve the interval $\rho(v)$ solely from the permutation, because we do not know one of the endpoints of the interval. Then for every node v in the sequence S we store the interval $\rho(v)$ explicitly.

Finally, after a reversal of an interval in π we can obtain at most two new adjacencies in π , at both ends of the interval, because pairs of positions fully inside or outside the interval cannot become or stop being an adjacency. Hence we can query elements adjacent to the endpoints of the reversed interval, check if an adjacency appeared and, if yes, remove the corresponding element from V . Hence, we can implement operations 1-4 from Theorem 2.6 using operations 1-4 from Theorem 2.7.

When the algorithm terminates, $OV(\pi)$ consists only of white isolated nodes. Furthermore, by construction, all chosen nodes were safe (otherwise we would have obtained an all-white connected component), so due to Theorem 2.5 the sequence is the shortest possible. \square

3 Improved algorithm for sorting signed permutations

In this section, we show how to efficiently implement all the four operations required in the interface presented in Theorem 2.7. Recall that π is a signed permutation on n elements and augment it with 0 and $n + 1$ at the ends that do not take part in any reversal.

Fact 3.1 ([12], Theorem 5). *There exists algorithm implementing operations 1 (query for π_i or π_i^{-1}) under reversals on π using $\mathcal{O}(\log n)$ time per query or reversal.*

Now we explain how to implement operations 2-4. We create the following graph $G(\pi)$ on $\mathcal{O}(n)$ nodes, $\mathcal{O}(n)$ blue edges and $n + 1$ red edges. There are two groups of nodes: $0, 1, 2, 3, \dots, (n + 1)$ and $0', 1', \dots, (n + 1)'$, where we think that both i and i' correspond to the element i of the permutation.

For every $0 \leq i \leq n$, we create a blue edge $b_i = \{i, (i + 1)\}$ and a blue edge $b'_i = \{i', (i + 1)'\}$. Additionally, if i and $i + 1$ have the same sign in π we create two red edges $r_i = \{i, (i + 1)\}$ and $r'_i = \{i', (i + 1)'\}$. Otherwise, if i and $i + 1$ have different sign in π we create two red edges $r_i = \{i, (i + 1)'\}$ and $r'_i = \{i', (i + 1)\}$. See Figure 3. Intuitively, we will use the blue edges to simulate reversals on the permutation (operation 3), and the red edges correspond to the elements of V (operations 2 and 4).

A *blue connected component* is a connected component of a subgraph of $G(\pi)$ containing all nodes of $G(\pi)$, but only the blue edges.

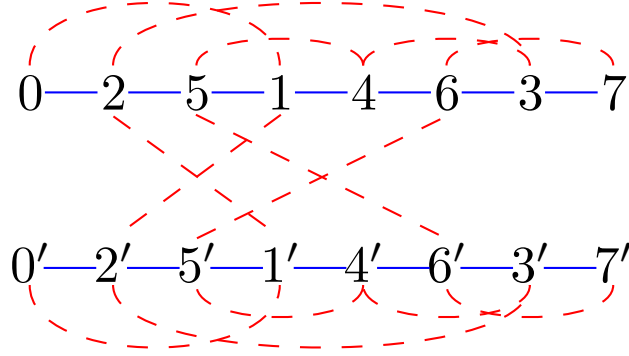


Figure 3: $G(\pi)$ for $\pi = (0, -2, -5, 1, -4, 6, -3, 7)$. Red edges are dashed and blue edges are solid.

Let c be the complementation function that satisfies $c(i) = i'$ and $c(i') = i$, for $i \in \{0, \dots, (n + 1)\}$. Our algorithm maintains the following invariants:

1. There are two blue connected components.
2. There is a blue edge $\{a, b\}$ if and only if there is a blue edge $\{c(a), c(b)\}$.
3. Every blue connected component is a path which, when being read from one of the endpoints (0 or $0'$), consists of nodes corresponding to $|\pi_0|, |\pi_1|, \dots, |\pi_{n+1}|$ in this order.

Clearly, the invariants hold before we apply any reversal. To simulate operation 3 from Theorem 2.7, that is, a reversal of an interval $[a, b]$, we perform an operation hereafter referred to as *reconnection* as follows. We reconnect the four blue edges in such a way that we reverse the nodes at the positions $a, a + 1, \dots, b$ on both paths and “change their sides”. See Figure 4 for an example. Formally, if $b_{a-1} = \{x, y\}$ and $b'_{a-1} = \{u, v\}$, where both y and v correspond to $|\pi(a)|$, we replace them with $b_{a-1} = \{x, v\}$ and $b'_{a-1} = \{u, y\}$. Analogously, if $b_b = \{p, t\}$ and $b'_b = \{s, q\}$, where p and s correspond to $|\pi(b)|$, we replace them with $b_b = \{p, q\}$ and $b'_b = \{s, t\}$. All other edges remain as they were. See Figure 5. Clearly, the above invariants are preserved after a reconnection. We note that translating a reversal into insertion and deletion of edges of $G(\pi)$ can be done efficiently using operation 1.

To simulate operation 4 (removal of an element i from V), we remove r_i and r'_i from $G(\pi)$. This operation does not violate the invariants either. We finally explain how to simulate operation 2 (find $i \in V$ such that i and $i + 1$ have different signs in π). Recall that no reversal includes element 0. Let $\#(i)$ be the number of performed reversals that changed the sign of i . The blue components fully describe the signs of the elements of π :

Lemma 3.2. *Node i is in the same blue component as node 0 if and only if $\#(i)$ is even.*

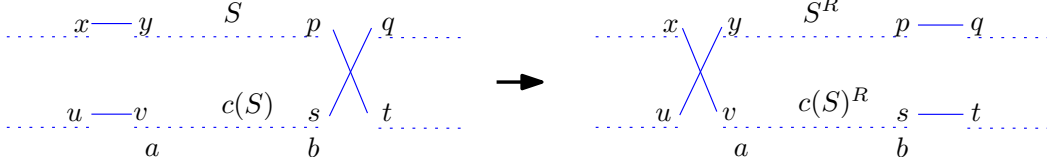


Figure 4: Reconnection of edges during the reversal of an interval $[a, b]$. By invariant 2, $c(x) = u, c(y) = v, c(p) = s$ and $c(q) = t$. $c(S)$ denotes complementing all nodes in S , and P^R denotes reversing the path P .

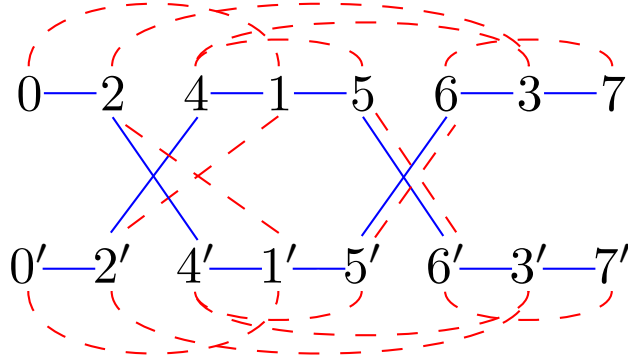


Figure 5: The result of reversing the interval $[2, 4]$ on the graph from Figure 3.

Proof. By induction on the number of performed reversals. By definition of $G(\pi)$, in the beginning every node i is in the same blue component as 0 and every node i' is in the same component as $0'$. Consider the graph after reversing an interval $[a, b]$. Observe that all nodes j and j' such that $\pi_j^{-1} < a$ or $\pi_j^{-1} > b$ are unaffected by the reversal (with regards to being in the same component as 0), so the claim follows for them by induction. On the other hand, nodes j and j' such that $a \leq \pi_j^{-1} \leq b$ change, that is after the reconnection they are in the component of 0 if and only if they were not before while the parity of $\#(j)$ changes. \square

We immediately derive a property that allows efficient implementation of operation 2 from Theorem 2.7:

Corollary 3.3. *After a sequence of reversals, the endpoints of a red edge r_i are in distinct blue components if and only if i and $i + 1$ have different signs in π . In this case, r_i is called active. Symmetrically for a red edge r'_i .*

Proof. By construction, the claim holds at initialisation. As the red edges are never affected by a reconnection, the claim follows by Lemma 3.2. \square

In other words, to implement operation 2 it suffices to find an active red edge in $G(\pi)$ under insertions and deletions of blue edges and deletion of red edges. As a warm-up, we show a very simple $\mathcal{O}(\log^4 n)$ amortized time algorithm.

Theorem 3.4. *There exists an algorithm supporting operations 2-4 from Theorem 2.7 on a permutation π on n elements in $\mathcal{O}(\log^4 n)$ amortized time.*

Proof. We assign weights to the edges of $G(\pi)$ in the following way: all blue edges get weight 0 and for every i , the edges r_i and r'_i get weight $i + 1$. Additionally, we add an edge between 0 and $0'$ with weight ∞ . As the 0-weight edges divide the graph into two components (recall the invariants of

the reconnection operation), with the additional edge $\{0, 0'\}$ the graph is always connected. Hence, its minimum spanning tree (MST) has total weight ∞ when there is no active edge. Otherwise, namely when there is at least one active red edge, its weight is w , where $w - 1$ is the smallest i such that r_i is active. Thus from the weight of the MST we can retrieve an active red edge. By applying an existing approach for dynamic MST that runs in $\mathcal{O}(\log^4 n)$ amortized time by Holm et al. [10], the claim follows. \square

Now we show a more efficient approach that does not use dynamic MST. First we remind some properties of link-cut trees.

Definition 3.5 (Link-cut tree [16]). *A link-cut tree is a data structure for maintaining a forest (a set of rooted trees) subject to the following operations:*

1. *make_tree(): Add a tree consisting of a single node to the forest.*
2. *cut(v): Cuts a tree containing a non-root node v into two by deleting the edge from the parent of v to v ;*
3. *link(v, w): Given two nodes v, w in different trees of the forest, links them by making v a child w . The operation assumes that v is the root of the tree it belongs to.*
4. *find_root(v): Finds the root of the tree containing a node v .*

Sleator and Tarjan [16] showed an implementation of link-cut trees with $O(\log n)$ amortized time per operation. Their implementation partitions each tree T in the forest into so-called *preferred paths*, defined as follows: A child w of a node v is the preferred child if the last access (to be described later) within v 's subtree was in w 's subtree. In particular, if the last access within v 's subtree was to v , it has no preferred child. A preferred edge is an edge between a preferred child and its parent. A preferred path is a maximal continuous path of preferred edges in a tree, or a single node if there is no preferred edge incident with it. On top of each preferred path, Sleator and Tarjan maintain a biased 2-3 search tree [3] (a balanced binary search tree), where the nodes are keyed by their depths in T .¹

All operations from Definition 3.5 are implemented with *access(v)*. If a node v belongs to a tree T of the forest, the subroutine simulates accessing the nodes in the path from the root of T to v by updating the preferred paths. As a corollary, if we access a node v and then a node w of T , the unique path between them will span at most two preferred paths, which will be important later.

Theorem 3.6. *There exists an algorithm supporting operations 2-4 from Theorem 2.7 on a permutation π on n elements in $\mathcal{O}(\log^2 n)$ amortized time.*

Proof. Consider the graph $G(\pi)$ again. On top of $G(\pi)$, we maintain a data structure for fully-dynamic graph connectivity [10], which supports edge insertions/deletions in $\mathcal{O}(\log^2 n)$ amortized time and maintains a spanning forest F of $G(\pi)$ in a link-cut tree.

There is an active red edge if and only if 0 and $0'$ are connected by a path, which we can check in $O(\log n)$ time using the link-cut tree. If 0 and $0'$ are connected, they belong to the same spanning tree $T \in F$, and are connected by a unique path $P = v_0 - v_1 - \dots - v_k$ of T , where $v_0 = 0$ and $v_k = 0'$. Since P starts in the blue component of 0 and ends in the blue component of $0'$, it must contain an active red edge, because 0 and $0'$ are in distinct blue components.

¹Modern variations of link-cut trees represent the paths with splay trees [17] to simplify the analysis, but we stick to the original biased 2-3 search trees which are guaranteed to have $\mathcal{O}(\log n)$ depth.

We binary search P for an active red edge using the following observation. Consider a fragment $v_i - v_{i+1} - \dots - v_j$ of P with the property that v_i belongs to the blue component of 0 while v_j belongs to the blue component of $0'$. Then, consider any $k \in \{i, i+1, \dots, j-1\}$. If v_k belongs to the blue component of $0'$ then we can recurse on $v_i - v_{i+1} - \dots - v_k$. If v_{k+1} belongs to the blue component of 0 then we can recurse on $v_{k+1} - v_{k+2} - \dots - v_j$. The remaining case is that $\{v_k, v_{k+1}\}$ is an active red edge.

To implement the search efficiently, we first call $access(0)$ and $access(0')$ to guarantee that P spans at most two preferred paths. In fact, as the preferred paths end at 0 or $0'$, P spans the path to 0 entirely, and spans a suffix of the path to $0'$. Thus, $P = v_0 - v_1 - \dots - v_m - v_{m+1} - \dots - v_k$, where $v_0 - v_1 - \dots - v_m$ is a fragment of one preferred path and $v_{m+1} - v_{m+2} - \dots - v_k$ is a fragment of the other preferred path. We first use the above observation to narrow down the search to one of them (or terminate after having found an active red edge). Then, we traverse in the biased 2-3 tree while using the above observation to decide if we should descend to the left or to the right child. In more detail, we first descend to the lowest common ancestor of the endpoints of the fragment. Then, we continue the descent, in every step using the above observation twice: if the current node is v_i , we apply it on $k = i - 1$ and $k = i$. We note that accessing v_{i-1} and v_{i+1} can be done in constant time by maintaining links to the predecessor/successor of each node on its path.

By Lemma 3.2, we can test if a node v_j is in the same blue component as node 0 by checking the sign of j in π , which in turn can be done in $\mathcal{O}(\log n)$ time by Fact 3.1. The biased 2-3 trees used to represent each path have depth $\mathcal{O}(\log n)$, so the search takes $\mathcal{O}(\log^2 n)$ time. \square

We are ready to present the final version of our algorithm. First, we observe that inside Theorem 3.6 we can use the faster fully dynamic connectivity structure, which supports edge insertions/deletions in $\mathcal{O}(\log^2 n / \log \log n)$ amortized time [20]. The structure can be extended to maintain a spanning forest F of $G(\pi)$ stored in a link-cut tree. Then, when searching for an active red edge we perform $\mathcal{O}(\log n)$ steps. If every step can be implemented in only $\mathcal{O}(\log n / \log \log n)$ time then the overall complexity becomes $\mathcal{O}(n \log^2 n / \log \log n)$. This is indeed possible thanks to the following lemma:

Lemma 3.7. *For any $\varepsilon > 0$, there is algorithm that maintains a signed permutation π on n elements under reversals in $\mathcal{O}(\log^{1+\varepsilon} n)$ time and, given an element i , retrieves $\#(i)$ in $\mathcal{O}(\log n / \log \log n)$ time.*

Proof. Our data structure is a slightly modified data structure from Fact 3.1 [12, Theorem 5] which we used to query for π_i or π_i^{-1} . We start by recalling its main components. It stores the elements of π in a balanced binary search tree allowing for splits and joins, for example a red-black tree. Additionally, at each node it stores the size of its subtree and the reverse flag, which, if set to true, indicates that the subtree rooted at that node should be read in reverse order, that is from right to left, and the signs of its elements should be flipped. It maintains the following invariant: the in-order traversal of the tree, taking the reverse flags into account, must give the permutation. To execute a reversal $\rho(i, j)$, we split the tree into three trees: T_1 containing π_1, \dots, π_{i-1} , T_2 containing π_i, \dots, π_j , and T_3 containing π_{j+1}, \dots, π_n . Next, we flip the reverse flag of the root of T_2 , and finally join T_1 and T_2 and then the resulting tree and T_3 . It is easy to see that the resulting tree satisfies the invariant. Via a standard approach, the update requires to perform $\mathcal{O}(\log n)$ rotations that touch $\mathcal{O}(\log n)$ nodes.

We augment the tree with shortcut pointers that slow down the updates but make the queries faster. For a node v of depth d , the shortcut points to its ancestor w of depth $\max\{0, d - \lfloor \varepsilon \log \log n \rfloor\}$ and stores the number of reverse flags on the path from v to w that are set to true. To compute $\#(i)$, we start from the node containing i and follow the shortcuts to the root aggregating the

number of the reverse flags that are set to true in $\mathcal{O}(\log n / \varepsilon \log \log n) = \mathcal{O}(\log n / \log \log n)$ time, as the depth of a node in a red-black tree on n elements is bounded by $\mathcal{O}(\log n)$. To perform a reversal, we use the algorithm described in the previous paragraph and then update the shortcuts for all nodes u of the tree such that one of the nodes between u and the endpoint of the shortcut for u was modified by a rotation.

We update the shortcuts in $\mathcal{O}(\log^{1+\varepsilon} n)$ total time in the following way. Consider a node x modified by a rotation. Let T_x be the subtree rooted at x . We must update the shortcuts for all nodes in T_x that have depth at most $\varepsilon \log \log n$. The number of such nodes is $\mathcal{O}(2^{\varepsilon \log \log n}) = \mathcal{O}(\log^\varepsilon n)$. The update consists of two steps. First, for each of the $\varepsilon \log \log n$ lowest ancestors of x , we compute the number of reverse flags that are set to true on the path from the ancestor to x in $\mathcal{O}(\varepsilon \log \log n)$ total time. Second, we run a depth-limited DFS from x that traverses all nodes of T_x that have depth at most $\varepsilon \log \log n$. During the DFS, we store the depth of the currently visited node v and the number of reverse flags that are set to true on the path from x to v , and update the shortcut from v by retrieving the appropriate ancestor of x . As the DFS takes time linear in the number of visited nodes, the total time of the two steps is $\mathcal{O}(\log^\varepsilon n)$. Hence, the update takes $\mathcal{O}(\log^{1+\varepsilon} n)$ across all $\mathcal{O}(\log n)$ nodes modified by a rotation. \square

By choosing $0 < \varepsilon < 1$ we obtain $\mathcal{O}(\log^{1+\varepsilon} n) = \mathcal{O}(\log^2 n / \log \log n)$ update time. Recall that in Theorem 3.6, operation 2 required $\mathcal{O}(\log n)$ queries about $\#(i)$, so applying the above data structure gives us the following result:

Corollary 3.8. *There exists an algorithm supporting operations 2-4 from Theorem 2.7 on a permutation π on n elements in $\mathcal{O}(\log^2 n / \log \log n)$ amortized time.*

To conclude, by combining the above corollary, Fact 3.1 and Theorem 2.7 we obtain an algorithm sorting a signed permutation by reversals in $\mathcal{O}(n \log^2 n / \log \log n)$ time.

References

- [1] David A. Bader, Bernard M. E. Moret, and Mi Yan. A linear-time algorithm for computing inversion distance between signed permutations with an experimental study. *J. Comput. Biol.*, 8(5):483–491, 2001.
- [2] Vineet Bafna and Pavel A. Pevzner. Genome rearrangements and sorting by reversals. *SIAM J. Comput.*, 25(2):272–289, 1996.
- [3] Samuel W. Bent, Daniel D. Sleator, and Robert E. Tarjan. Biased 2-3 trees. In *FOCS*, pages 248–254, 1980.
- [4] Anne Bergeron. A very elementary presentation of the Hannenhalli-Pevzner theory. *Discret. Appl. Math.*, 146(2):134–145, 2005.
- [5] Piotr Berman and Sridhar Hannenhalli. Fast sorting by reversal. In *CPM*, volume 1075 of *Lecture Notes in Computer Science*, pages 168–185. Springer, 1996.
- [6] Alberto Caprara. Sorting by reversals is difficult. In *RECOMB*, pages 75–83. ACM, 1997.
- [7] William H.E. Day and David Sankoff. Computational complexity of inferring phylogenies from chromosome inversion data. *Journal of Theoretical Biology*, 124(2):213–218, 1987.

- [8] Yijie Han. Improving the efficiency of sorting by reversals. In *BIOCOMP*, pages 406–409. CSREA Press, 2006.
- [9] Sridhar Hannenhalli and Pavel A. Pevzner. Transforming cabbage into turnip: Polynomial algorithm for sorting signed permutations by reversals. *J. ACM*, 46(1):1–27, 1999.
- [10] Jacob Holm, Kristian de Lichtenberg, and Mikkel Thorup. Poly-logarithmic deterministic fully-dynamic algorithms for connectivity, minimum spanning tree, 2-edge, and biconnectivity. *J. ACM*, 48(4):723–760, 2001.
- [11] Haim Kaplan, Ron Shamir, and Robert Endre Tarjan. A faster and simpler algorithm for sorting signed permutations by reversals. *SIAM J. Comput.*, 29(3):880–892, 1999.
- [12] Haim Kaplan and Elad Verbin. Sorting signed permutations by reversals, revisited. *J. Comput. Syst. Sci.*, 70(3):321–341, 2005.
- [13] John D. Kececioglu and David Sankoff. Exact and approximation algorithms for sorting by reversals, with application to genome rearrangement. *Algorithmica*, 13(1/2):180–210, 1995.
- [14] Michal Ozery-Flato and Ron Shamir. Two notes on genome rearrangement. *J. Bioinform. Comput. Biol.*, 1(1):71–94, 2003.
- [15] Pavel A. Pevzner. *Computational molecular biology - an algorithmic approach*. MIT Press, 2000.
- [16] Daniel D. Sleator and Robert Endre Tarjan. A data structure for dynamic trees. In *STOC*, page 114–122, 1981.
- [17] Daniel Dominic Sleator and Robert Endre Tarjan. Self-adjusting binary search trees. *J. ACM*, 32(3):652–686, jul 1985.
- [18] Krister M. Swenson, Vaibhav Rajan, Yu Lin, and Bernard M. E. Moret. Sorting signed permutations by inversions in $O(n \log n)$ time. *J. Comput. Biol.*, 17(3):489–501, 2010.
- [19] Eric Tannier, Anne Bergeron, and Marie-France Sagot. Advances on sorting by reversals. *Discret. Appl. Math.*, 155(6-7):881–888, 2007.
- [20] Christian Wulff-Nilsen. Faster deterministic fully-dynamic graph connectivity. In *SODA*, page 1757–1769. SIAM, 2013.

A Proof of Theorem 2.6

The goal of this section is to establish correctness of Algorithm 1. We will first analyse its slower version, see Algorithm 2. Then, we will consider the faster version, see Algorithm 3, and show that the two algorithms are in fact equivalent. Finally, we will notice that Algorithm 3 is a reformulation of Algorithm 1 that avoids the goto statement.

Consider the following algorithm. We will first show that, given a graph with no non-singleton all-white components, it construct a sequence of toggles that makes all nodes white and removes all edges from the graph. We note that with every node v in any sequence we store the neighborhood of v in the graph before toggling v . This allows us to undo `toggle(v)` in the future.

The overall structure of the proof is as follows. We start with some technical lemmas concerning the properties of the algorithm. Then, we establish that in line 9 such a split indeed exists, and

Algorithm 2

```
1: function PROCESS(graph  $G$  with no non-singleton all-white connected components)
2:    $S := ()$ 
3:    $G' := G$ 
4:   while there is a black node  $v$  in  $G'$  do
5:     apply toggle( $v$ ) to  $G'$ 
6:      $S := S, v$ 
7:   while there is a non-singleton all-white connected component in  $G/S$  do
8:      $U :=$  set of nodes from non-singleton all-white connected components in  $G/S$ 
9:      $S_1, S_2 := S$  for the longest  $S_1$  s.t.  $G/S_1$  has a node of  $U$  in a not-all-white component
10:     $G_1 := G/S_1$ 
11:     $S_3 := ()$ 
12:    while there is a black node  $v$  from  $U$  in  $G_1$  do
13:      apply toggle( $v$ ) to  $G_1$ 
14:       $S_3 := S_3, v$ 
15:    if  $S_2[1]$  is white in  $G_1$  then
16:      remove the last element  $w$  from  $S_3$ 
17:      undo toggle( $w$ ) in  $G_1$ 
18:     $S := S_1, S_3, S_2$ 
19:  return  $S$ 
```

further $v_{mid} = S_2[1]$ does exist. Next, we show that in line 16 the sequence S_3 is non-empty. In fact, its length is always at least 2, which immediately gives us that the whole algorithm terminates in a finite number of steps (in fact, at most linear in the number of nodes). To prove the correctness, we call a sequence S of toggles *valid* for a graph G if we can toggle nodes of G in the order defined by S , that is s_i is a black node in $G/s_1/\dots/s_{i-1}$ for every $1 \leq i \leq k$. First, we will show that in line 7 S is valid for G . Next, we will argue that G/S does not contain any black nodes. Finally, we will prove that the algorithm terminates in a finite number of steps (in fact, at most linear in the number of nodes). By condition in line 7, this implies that the found sequence of toggles makes all nodes white and isolated as required.

Within a single iteration of the main while loop in lines 7-18 of Algorithm 2, let $v_{mid} = S_2[1]$, O be the set of nodes in not-all-white components of G_1/v_{mid} , and $L = \Gamma(v_{mid}) \cap O$.

Claim A.1. *Before line 12 of Algorithm 2, U consists of all nodes in the non-singleton all-white connected components of G_1/v_{mid} .*

Proof. Consider a node v in G_1/v_{mid} . If v is white and isolated in G_1/v_{mid} , it will be also white and isolated in G/S , so it does not belong to U . If v belongs to a white non-singleton component of G_1/v_{mid} , its component cannot be affected by any further toggle and remains white in G/S , so $v \in U$. Finally, $v \in O$ (a node that belongs to a not-all-white component) cannot belong to U as in this case S_1, v_{mid} satisfies the conditions of the split and is longer than S_1 , see line 9. \square

Consequently, the nodes of G_1 are v_{mid} , some white isolated nodes, and the nodes from O and U . Next, we show the following properties of G_1 that are preserved after any sequence of toggles of black nodes from U :

Lemma A.2. *Consider the graph G_1 before line 12 of Algorithm 2, and any sequence w_1, w_2, \dots, w_k of nodes from U such that for $1 \leq i \leq k$, w_i is a black node in $G_1/w_1/w_2/\dots/w_{i-1}$. The graph $G_1/w_1/w_2/\dots/w_k$ satisfies each of the following properties:*

- (1) A node of U is black if and only if it is adjacent to v_{mid} .
- (2) There are all possible edges between $\Gamma(v_{mid}) \cap U$ and L .
- (3a) There is no edge between the nodes of $U \setminus \Gamma(v_{mid})$ and the nodes outside U .
- (3b) There is no edge between the nodes of $O \setminus L$ and the nodes outside O .

The sets O, U and L are defined for G_1 , whereas adjacency and $\Gamma(v_{mid})$ refer to $G_1/w_1/w_2/\dots/w_k$.

Proof. Induction on k , following from the properties of the toggle operation.

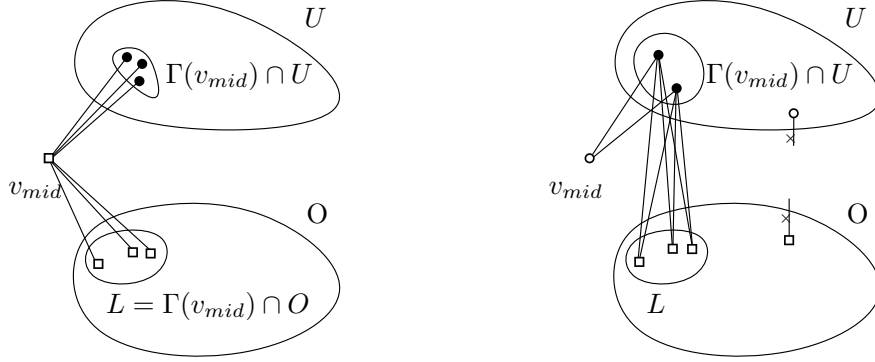


Figure 6: Full nodes denote black nodes, empty nodes denote white nodes, and square nodes denote nodes which can be either black or white. O is the set of nodes in the not-all-white components of G_1/v_{mid} , and U is the set of nodes from the non-singleton all-white connected components in G/S . Left: v_{mid} and the sets O, L, U . Right: invariants (1)-(3b). \times denotes the edges that are not present in the graph.

Base case ($k = 0$). We prove each of the properties separately:

- (1) Otherwise in G_1/v_{mid} we would have a black node from U , which contradicts Claim A.1.
- (2) Otherwise in G_1/v_{mid} we would have an edge from a node $u \in U$ to a node from O , which means that u is in a not-all-white component of G_1/v_{mid} , which contradicts Claim A.1.
- (3a) If in G_1 there was an edge from a node $u \in U \setminus \Gamma(v_{mid})$ to a node outside U , in G_1/v_{mid} this edge will be preserved and u will be part of a not-all-white component, which contradicts Claim A.1.
- (3b) If in G_1 there was an edge from a node $t \in O \setminus L$ to a node u outside O , then $u \in U$ as $t \notin \Gamma(v_{mid})$ and hence $u \neq v_{mid}$. In G_1/v_{mid} this edge will be preserved and u will be part of a not-all-white component, which contradicts Claim A.1.

Induction step ($k - 1 \rightarrow k$). Let $G' = G_1/w_1/w_2/\dots/w_{k-1}$ and $G'' = G'/w_k$. By definition of w_k and (1) for G' , the node w_k is a neighbor of v_{mid} in G' . By (2) for G' , the node w_k is connected to all nodes from L in G' .

- (1) Nodes not adjacent to w_k in G' are not affected by $\text{toggle}(w_k)$ so the claim follows for them by induction. Consider a neighbor v of w_k . If $v \in \Gamma(v_{mid})$ then by (1) v is black. Hence v is white in G'' and $v \notin \Gamma(v_{mid})$. If $v \notin \Gamma(v_{mid})$ then by (1) v is white. Hence v is black in G'' and $v \in \Gamma(v_{mid})$.

(2) There are two cases in which $v \in \Gamma(v_{mid})$ in $G'' = G'/v_{mid}$:

- $v \in \Gamma(w_k), v \notin \Gamma(v_{mid})$: by (3a), in G' the node v is connected to none of the nodes from O . As w_k is connected to v_{mid} and all nodes from L , in G'' the node $v \in \Gamma(v_{mid})$ and is connected to all nodes from L , so (2) holds.
- $v \notin \Gamma(w_k), v \in \Gamma(v_{mid})$: v is not affected by $\text{toggle}(w_k)$ in G' , so $v \in \Gamma(v_{mid})$ in G'' and (2) follows by induction.

(3a) There are two cases in which $v \notin \Gamma(v_{mid})$ in $G'' = G'/v_{mid}$:

- $v \in \Gamma(w_k), v \in \Gamma(v_{mid})$: in G' the node v is connected to all nodes from L by (2), but no node from $O \setminus L$ by (3b). Hence $v \notin \Gamma(v_{mid})$ in G'' and v is not connected to any node from O , so (3a) holds.
- $v \notin \Gamma(w_k), v \notin \Gamma(v_{mid})$: v is not affected by $\text{toggle}(w_k)$, so in G'' the node $v \notin \Gamma(v_{mid})$ and (3a) follows by induction.

(3b) By (3b) for G' , nodes from $O \setminus L$ are not neighbors of w_k in G' , so they are not affected by $\text{toggle}(w_k)$ and (3b) holds also for G'' . \square

Now note that after passing the check in line 7 of Algorithm 2, the set U is non-empty. In line 9, we can always find such a split, because in particular $S_1 = ()$, $S_2 = S$ is a valid split, as G does not have non-singleton white components. Furthermore, by definition of the split, in line 15 we have $S_2 \neq ()$, so $v_{mid} = S_2[1]$ does exist. Moreover, by maximality of S_1 :

Claim A.3. $S_2[2\dots]$ contains only nodes from O .

Proof. Assume otherwise, and let $S[i]$ be the last node in $S_2[2\dots]$ that does not belong to O . Then it must belong to U , as every other node is white already in G_1/v_{mid} . Thus, we could have split S into $S_1, v_{mid}, S_2[2\dots(i-1)]$ and $S[i\dots]$. \square

Now we show that Algorithm 2 is well-defined and terminates.

Lemma A.4. After the while loop in lines 12-14 of Algorithm 2, S_3 contains at least 2 nodes.

Proof. If $|S_3| = 0$ and there are no black nodes from U in G_1 , by Lemma A.2(1) there are no neighbors of v_{mid} in U and by (3) there are no edges from U to O . Therefore, the connected components of U are not affected by $\text{toggle}(G_1, v_{mid})$ and G_1 has a node of U in a not-all-white component if and only if G_1/v_{mid} has a node of U in a not-all-white component, which contradicts the maximality of S_1 in line 9.

If $|S_3| = 1$, there is a black node $w \in U$ such that G_1/w has no black nodes in U . By Lemma A.2(1), the node $w \in \Gamma(v_{mid})$ and $\Gamma^+(w) \cap U = \Gamma(v_{mid}) \cap U$ as there are no black nodes in U in G_1/w . Next, by Lemma A.2(2), w is adjacent to all nodes from L and by Lemma A.2(3b) adjacent to no nodes from $O \setminus L$. Therefore, in G_1/v_{mid} the node w is white and isolated which is the case also for $G_1/S_2 = G/S$. Hence $w \notin U$, a contradiction. \square

In particular, in line 16 we have $S_3 \neq ()$, so we can indeed remove the last element from S_3 . Further, in line 18 we have $|S_1, S_3, S_2| > |S|$, so each iteration of the main while loop results in a longer S . This implies that the algorithm terminates, as each element of S makes at least one node of G/S isolated.

To show the next property, let $W = L \cup \{v_{mid}\}$ and let $G|_Y$ denote the subgraph of graph G induced by a set Y of nodes.

Claim A.5. Before line 15 of Algorithm 2, the subgraph $(G_1/S_3)|_W$ equals $G_1|_W$ if $|S_3|$ is even. If $|S_3|$ is odd, the subgraph of $(G_1/S_3)|_W$ is the local complement of $G_1|_W$.

Proof. By induction on $|S_3|$. By Lemma A.2, every $w_i \in S_3$ is adjacent to v_{mid} and all nodes from L . Therefore, every operation $\text{toggle}(w_i)$ complements the subgraph of the current graph induced by W and the claim follows. \square

We note that before we remove the last node w_{last} in line 16 of the algorithm, w_{last} has exactly the same neighborhood as v_{mid} , so removing it and next toggling v_{mid} gives exactly the same result:

Claim A.6. Consider the graph G_1 before line 12 of Algorithm 2, and let S_3, w_{last} be a valid sequence of nodes from U for G_1 such that there is no black node from U in $G_1/S_3/w_{last}$. If $|S_3|$ is even, in G_1/S_3 we have $\Gamma^+(w_{last}) = \Gamma^+(v_{mid})$. If $|S_3|$ is odd, v_{mid} has no neighbors from U in $G_1/S_3/w_{last}$.

Proof. By definition, there is no black node from U in $G_1/S_3/w_{last}$, so by Lemma A.2(1) v_{mid} has no neighbors from U in $G_1/S_3/w_{last}$ and the claim follows for odd $|S_3|$.

If $|S_3|$ is even, by Claim A.5, v_{mid} is adjacent to all nodes from L and is black (because it is black in G_1). By Lemma A.2, $\Gamma^+(w_{last}) \cap O = \Gamma^+(v_{mid}) \cap O = L$. Recall that w_{last} is adjacent to v_{mid} . Consider a node $v \in U$. If $v \in \Gamma^+(w_{last})$, then it is black in G_1/S_3 , because in $G_1/S_3/w_{last}$ there is no black node from U . Therefore, $v \in \Gamma^+(v_{mid})$, by Lemma A.2(1). If $v \notin \Gamma^+(w_{last})$, then it is white in G_1/S_3 , because in $G_1/S_3/w_{last}$ there is no black node from U , so $v \notin \Gamma^+(v_{mid})$, by Lemma A.2(1). \square

Lemma A.7. In line 18 of Algorithm 2, S_1, S_3, S_2 is a valid sequence for G .

Proof. By definition of the split $S = S_1, S_2$ and Lemma A.2, toggling the nodes $S_2[1 \dots]$ affects only the subgraph of $(G_1/v_{mid})|_O$, because there are no edges from O to outside O and all other nodes are either isolated or in all-white connected components.

Clearly, S_3 is a valid sequence for G_1 . As v_{mid} is black in G_1 (toggling it is the first operation in S_2), by Claim A.5 the subgraph $(G_1/S_3)|_W = G_1|_W$ if and only if v_{mid} is black in G_1/S_3 . In line 15 of Algorithm 2, we ensure that v_{mid} is black in G_1/S_3 , so then S_3, v_{mid} is a valid sequence for G_1 . As the subgraph $G_1|_{O \setminus L}$ is not affected by the operations from S_3 , we obtain that $(G_1/S_3)|_{O \cup \{v_{mid}\}} = G_1|_{O \cup \{v_{mid}\}}$.

To sum up, the subgraph $(G_1/v_{mid})|_O$ is the same as the subgraph $(G_1/S_3/v_{mid})|_O$. As toggling the nodes $S_2[1 \dots]$ affects only the subgraph $(G_1/v_{mid})|_O$, we obtain that S_3, S_2 is a valid sequence for G_1 and the claim follows. \square

Let S'_3 be the sequence S_3 when reaching line 15 and S'' be the sequence S_3 when reaching line 18.

Lemma A.8. Consider the graph G_1 before line 12 of Algorithm 2. The following properties hold:

1. $G_1/v_{mid}|_O = G_1/S''_3/v_{mid}|_O$.
2. In $G_1/S''_3/v_{mid}$, the nodes from $S_2[2 \dots]$ do not have any neighbors in U .
3. $G_1/S'_3|_U = G_1/S''_3/S_2|_U$.

Proof. From Lemma A.2 every toggled node $v \in S''_3$ is incident to all nodes from L and no nodes from $O \setminus L$. By Claims A.5 and A.6, property (1) follows.

By Lemma A.2, in G_1/S'_3 there is no edge from U to O , so by Claim A.6 this is also the case for $G_1/S''_3/v_{mid}$. As every element of $S_2[2 \dots]$ is in O by Claim A.3, property (2) follows. Finally, property (3) follows immediately from property (2) and Claim A.6. \square

This immediately gives us that in line 18 of Algorithm 2, the sequence $S_3, v_{mid}, S_2[2\dots]$ is valid for G_1 , so S_1, S_3, S_2 is valid for G . Thus by induction S is valid for G in line 7 of Algorithm 2. Further, we can show by induction that there is no black node in G/S .

Lemma A.9. *There are no black nodes in G/S .*

Proof. The proof is by induction on the number of iterations of the main loop in lines 7-18 of Algorithm 2. Before line 7 the claim is true because we have kept toggling black nodes as long as possible.

Consider a single iteration of the main loop and graph G_1 before line 12. Let $G_2 = G_1/S_3''/S_2$. After line 14, there is no black node from U in G_1/S_3' , so by Lemma A.8(3), there is no black node from U in G_2 . Clearly v_{mid} is white in G_2 , because $v_{mid} = S_2[1]$ is toggled. Now we need to show that there is no black node from O in G_2 . By Lemma A.8(1) and the fact that $S_2[2\dots]$ contains only nodes from O (Claim A.3) we have $G_2|_O = G_1/S_3''/v_{mid}/S_2[2\dots]|_O = G_1/v_{mid}/S_2[2\dots]|_O = G_1/S_2|_O = G/S|_O$ which by induction has no black node. Hence G_2 has no black node in O and the claim follows. \square

This concludes the proof of the correctness of Algorithm 2. The bottleneck of Algorithm 2 is that in every iteration of the loop in lines 7-18, we need to find the split $S = S_1, S_2$. To implement Algorithm 2 efficiently, we first establish that the sets U and the suffixes S_2 are monotone in subsequent iterations of the algorithm. Let $S^{(j)}$ and $U^{(j)}$ denote the sequence S and the set U at the beginning of the j -th iteration of the main loop of Algorithm 2. We first show the following:

Lemma A.10. $U^{(j+1)} = U^{(j)} \setminus \{v : v \text{ became white and isolated while toggling nodes in line 13 of } j\text{-th iteration of the main loop of Algorithm 2}\}$.

Proof. By Lemma A.8(1,2) and Claim A.6, graphs $G/S^{(j)}$ and $G/S^{(j+1)}$ differ only on some nodes from $U^{(j)}$, that is only nodes from $U^{(j)}$ can have different colors and both graphs have exactly the same edges except possibly of the edges that have both endpoints in $U^{(j)}$. Consider a node $v \in U^{(j)}$ toggled in line 13. All nodes that became white and isolated in G_1 while toggling v are white and isolated in $G/S^{(j+1)}$, so they do not belong to $U^{(j+1)}$ and the claim follows. \square

Similarly, let $S_2^{(j)}$ denote the sequence S_2 in the j -th iteration of the main loop of Algorithm 2. We show:

Lemma A.11. $S_2^{(j)}$ is a (not necessarily proper) suffix of $S_2^{(j+1)}$.

Proof. Consider the j -th iteration of the main while loop and let S_3' be the sequence S_3 when reaching line 15, S_3'' be the sequence S_3 when reaching line 18. Let $U^{(j)}$ and $U^{(j+1)}$ be defined as in Lemma A.10.

By the condition in line 12, there are no black nodes from $U^{(j)}$ in G_1/S_3' , and by Lemma A.2 there are no edges between the nodes of $U \setminus \Gamma(v_{mid})$ and the nodes outside $U^{(j)}$. Then, by Claim A.6, there are no black nodes from $U^{(j)}$ in $G_1/S_3''/v_{mid}$, and further there are no edges between the nodes of $U^{(j)}$ and the nodes outside $U^{(j)}$ there. Thus, no further toggles can affect the nodes from $U^{(j)}$. We conclude that, at the end of the j -th iteration, for every $i = 2, 3, \dots$, there are no black nodes in the connected component containing the nodes from $U^{(j)}$ in $G_1/S_3''/v_{mid}/S_2^{(j)}[2\dots i]$. However, by the previous lemma $U^{(j+1)} \subseteq U^{(j)}$, so this in particular holds for the nodes from $U^{(j+1)}$. Hence, the split in the j -th iteration must be before $v_{mid}, S_2^{(j)}[2\dots i] = S_2^{(j)}$. \square

This allows us to keep track of sequence S split into S_1 and S_2 and either move elements from the end of S_1 to the beginning of S_2 or add nodes at the end of S_1 . With this observation, we transform Algorithm 2 to Algorithm 3. Its correctness follows from the correspondence between appropriate variables in subsequent iterations of the main while loops:

Lemma A.12. *After the loop in lines 17-19 in the j -th iteration of the main loop of Algorithm 3 and after lines 8-11 in the j -th iteration of the main loop of Algorithm 2 the following are equal:*

- sequence S_1 in Algorithm 3 and sequence S_1 in Algorithm 2,
- sequence S_2 in Algorithm 3 and sequence S_2 in Algorithm 2,
- graph G in Algorithm 3 and graph G_1 in Algorithm 2,
- set V in Algorithm 3 and set U in Algorithm 2,

After the lines 20-23 in the j -th iteration of the main loop of Algorithm 3 and after lines 12-17 in the j -th iteration of the main loop of Algorithm 2 the following are equal:

- sequence S_1 in Algorithm 3 and sequence S_1, S_3 in Algorithm 2,
- sequence S_2 in Algorithm 3 and sequence S_2 in Algorithm 2,
- graph G in Algorithm 3 and graph G_1 in Algorithm 2,
- set V in Algorithm 3 and set U in Algorithm 2,

Proof. We proceed by induction on the number of iterations. The invariants on S_1, S_2, S_3 hold by the monotonicity shown in Lemma A.11. The invariant on G and G_1 holds trivially, because we always do or undo the corresponding toggles. Finally, to relate V and U we observe that the invariant initially holds, because we have removed from V all nodes that became isolated and white in lines 14-15 of Algorithm 3, and there are no black nodes in G/S_1 by condition of the while loop. Then, recall that we have already characterised how the set U changes in subsequent iterations in Lemma A.10, and this is exactly how V is being updated. \square

Finally, rearranging the order of while loops and checks gives us exactly Algorithm 1 from [19].

Algorithm 3

```
1: function PROCESS(graph  $G$  with no non-singleton all-white connected components)
2:    $S_1, S_2 := ()$ 
3:    $V :=$  the set of all nodes of  $G$  that are non-isolated or black

4:   function TOGGLEBLACKNODEFROMVANDAPPENDTOS1()
5:      $v :=$  a black node from  $V$ 
6:     toggle( $v$ )
7:     remove from  $V$  all nodes that became isolated and white (in particular: node  $v$ )
8:      $S_1 := S_1, v$ 

9:   function UNDOANDREMOVELASTMOVEFROMS1()
10:     $w :=$  last element of  $S_1$ 
11:    remove last element (that is:  $w$ ) from  $S_1$ 
12:    undo toggle( $w$ )
13:    return  $w$ 

14:   while there is a black node in  $V$  do
15:     TOGGLEBLACKNODEFROMVANDAPPENDTOS1()
16:   while  $V$  is non-empty do
17:     while there is no black node in  $V$  do
18:        $w :=$  UNDOANDREMOVELASTMOVEFROMS1()
19:        $S_2 := w, S_2$ 
20:     while there is a black node in  $V$  do:
21:       TOGGLEBLACKNODEFROMVANDAPPENDTOS1()
22:     if  $S_2[1]$  is white then
23:       UNDOANDREMOVELASTMOVEFROMS1()
24:   return  $S_1, S_2$ 
```
