

Prefix transpositions on binary and ternary strings

Amit Kumar Dutta, Masud Hasan, M. Sohel Rahman

Department of Computer Science & Engineering, Bangladesh University of Engineering and Technology, Dhaka, Bangladesh

Abstract: The problem of Sorting by Prefix Transpositions asks for the minimum number of prefix transpositions required to sort the elements of a given permutation. In this paper, we study a variant of this problem where the prefix transpositions act not on permutations but on strings over an alphabet of fixed size. Here, we determine the minimum number of prefix transpositions required to sort the binary and ternary strings, with polynomial time algorithms for these sorting problems.

Keywords: prefix transposition, ternary strings, binary strings, genome rearrangement

1 Introduction

The *transposition distance* between two permutations (and the related problem of *sorting by transposition*) is used to estimate the number of global mutations between genomes and can be used by molecular biologists to infer evolutionary and functional relationships. A transposition involves swapping two adjacent substrings of the permutation. In a prefix transposition, one of them must be a prefix. *Sorting by prefix transposition* is the problem of finding the minimum number of prefix transpositions needed to transform a given permutation into the identity permutation. In the literature, other interesting problems include sorting by other operations like reversals, prefix reversals, block interchange etc.

A natural variant of the aforementioned sorting problems is to consider them not on permutations but on strings over fixed size alphabets. This shift is inspired by the biological observation that multiple “copies” of the same gene can appear at various places along the genome [4]. Indeed, recent works by Christie and Irving [2], Radcliffe et al. [5] and Hurkens et al. [4] explore the consequences of switching from permutations to strings. Notably, such rearrangement operations on the strings have been found to be interesting and important in the study of orthologous gene assignment [1], especially if the strings have only low level of symbol repetition.

Chen et al. [1], presented for both reversals and transpositions, polynomial-time algorithms for computing the minimum number of operations to sort a given binary string. They also gave exact constructive diameter results on binary strings. Radcliff et al. [5] on the other hand gave refined and generalized reversal diameter results for non fixed size alphabets. Additionally, they gave a polynomial-time algorithm for optimally sorting a ternary (3 letter alphabet) string with reversals. Finally, Hurkens et al. [4] introduced *grouping* (a weaker form of sorting), where identical symbols need only be grouped together, while a group can be in any order. In the sequel, they gave a complete characterization of the minimum number of prefix reversals required to group (and sort) binary and ternary strings.

In this paper, we follow up the work of [4] and consider prefix transposition to group and sort binary and

ternary strings. Notably, as a future work in [4], the authors raised the issue of considering other genome arrangement operators. In particular, here, we find the minimum number of prefix transpositions required to group and sort binary or ternary strings. It may be noted that, apart from being an useful aid for sorting, grouping itself is a problem of interest in its own right [3].

The rest of the paper is organized as follows. In Section 2, we discuss the preliminary concepts and discuss some notations we use. Section 3 is devoted to grouping, where we present and prove the corresponding bounds. In Section 4, we present an algorithm to group the strings and in Section 5, we classify the strings along two different dimensions. Section 6 extends the results on grouping to get the corresponding bounds for sorting. Finally, we briefly conclude in Section 7.

2 Preliminaries

We follow the notations and definitions used in [4], which are briefly reviewed below for the sake of completeness. We use $[k]$ to denote the first k non-negative integers $\{0, 1, \dots, k-1\}$. A k -ary string is a string over the alphabet $\Sigma = [k]$. Moreover, a string $s = s_1s_2\dots s_n$ of length n is said to be *fully k -ary*, or to have *arity k* , if the set of symbols occurring in it is $[k]$.

A prefix transposition $f(1, x, y)$ on string s of length n , where $1 < x < y \leq (n+1)$, is an rearrangement event that transforms s into $[s[x] \dots s[y-1]s[1] \dots s[x-1]s[y] \dots s[n]]$. The prefix transposition distance $d_s(s)$ of s is defined as the number of prefix transpositions required to sort the string. Note that, after a transposition operation is performed, the two adjacent symbols of the corresponding string may be identical. We consider two strings to be equivalent if one can be transformed into the other by repeatedly duplicating (by transposing) symbols and eliminating adjacent identical symbols. This elimination of adjacent identical symbols gives us a *reduced string*, i.e., a string of reduced length and this process is referred to as *reduction*. As representatives of the *equivalence classes* we take the shortest string of each class. Clearly, these are the strings where adjacent symbols always differ. The process of transforming a string into the representative string of its equivalence class is sometimes referred to as *normalization*. Therefore, the

process of normalization basically comprises of repeated transposition and reduction.

For example, let $s = bababab$ and we want to apply operation $f(1, 3, 6)$. Now, $s[x] \dots s[y-1] = s[3] \dots s[5] = bab$, $s[1] \dots s[x-1] = s[1] \dots s[2] = ba$, $s[y] \dots s[n] = s[6] \dots s[7] = ab$. Therefore, after applying the operation, we get, $s = s[3] \dots s[5]s[1] \dots s[2]s[6] \dots s[7] = \overline{babbaab} = ba \mathbf{bbaa} b = ba \mathbf{ba} b = babab$

A reduction that decreases the string length by x is called an x -transposition. So, if $x = 0$, then we have a 0 transposition. The above example illustrates a 2-transposition.

3 Grouping

The task of sorting a string can be divided into two subproblems, namely, *grouping* the identical symbols together and then putting the groups of identical symbols in the right order. The grouping distance $d_g(s)$ of a fully k ary string s is defined as the minimum number of prefix transposition required to reduce the string to one of length k .

3.1 Grouping Binary Strings

As strings are normalized, only 2 kinds of binary strings are possible, namely, 010101...010 and 101010...101. The grouping of binary strings seems to be quite easy and obvious. The following bound is easily achieved.

Theorem 1. (Bound for Binary strings) *Let s be a fully binary string. Then, $d_g(s) = \lceil \frac{n-k}{2} \rceil$.*

Proof. We can always have a 2-transposition if $|s|$ is even. However, if $|s|$ is odd, we need an extra 1-transposition. So, the upper bound is $d_g(s) = \lceil \frac{n-k}{2} \rceil$. \square

We illustrate the above result with the help of an example. Let $s = ababab$. Then we can continue as follows: $s = ababab = \overline{ababab} = \mathbf{aabbab} = abab = \overline{abab} = \mathbf{aabb} = ab$. Here, we need two 2-transpositions to group this string. So, $d_g(s) = 2$.

3.2 Grouping Ternary Strings

In this section, we focus on ternary strings. As it seems, grouping ternary strings is not as easy as grouping binary strings. We start with the following theorem.

Lemma 1. *In a fully ternary string, we can always perform a 1-transposition.*

Proof. We take a ternary string s of length $n > 3$. Now, we take a prefix a of length k . If $a[1]$ occurs at the suffix at position i , we can transpose $a[1] \dots a[i-1]$ with $a[i]$. Then, $a[1]$ and $a[i]$ are adjacent and we can eliminate one of the two. Otherwise, if $a[k]$ occurs at the suffix at position i , then we can transpose $a[1] \dots a[k]$ with $a[k+1] \dots a[i-1]$. Then $a[k]$ and $a[i]$ are adjacent and as before, we can eliminate one of them. Since, one of the above cases always occurs for ternary strings, the result follows. \square

3.3 Grouping distance for Ternary strings

The lower bound for the *grouping* of a ternary string remains the same as that of binary strings; but, as can be seen from Theorem 2 below, the upper bound differs. We first give an easy but useful lemma.

Lemma 2. *Suppose $s[1..n]$ is a fully ternary string. If we have a prefix $s[1..i]$, $1 < i \leq n-2$ such that $s[1] = s[n-1]$ and $s[i] = s[n]$, then we have a 2-transposition.*

The proof of Lemma 2 is very easy and hence is omitted.

Theorem 2. (Bound for Ternary strings) *Let s be a fully ternary string. Then, $\lceil \frac{n-k}{2} \rceil \leq d_g(s) \leq \lceil \frac{n-k}{2} \rceil + 1$ where n is the length of the string and k is the arity.*

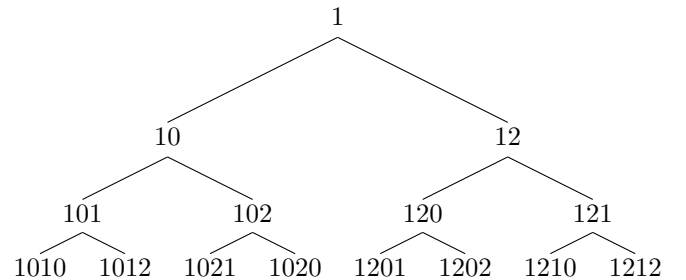
Proof. First we prove the lower bound. Here, $k = 3$. If we can always give a 2-transposition, then at each operation the string length is decremented by 2 and if $|s|$ is odd, we need an extra 1 transposition. Hence, we have $\lceil \frac{n-k}{2} \rceil \leq d_g(s)$.

Let us concentrate on the upper bound. As strings are fully ternary, we don't need to work with $n \leq 3$. Now, if we apply the upper bound for $n = 4, 5$ and 6 , we have $d_g(s) = 2, 2$ and 3 respectively. It is easy to realize that, by Lemma 1, we can always satisfy the above upper bound. Thus the upper bound is proved for $n < 7$.

Now we consider $n \geq 7$. In what follows, we only consider strings starting with 1. This doesn't lose the generality since we can always use relabeling for strings starting with 0 or 2. Now, note carefully that for any string starting with 1, we can only have one of the following eight prefixes of length 4:

$$1012, 1010, 1021, 1020, 1201, 1202, 1210 \text{ and } 1212. \quad (1)$$

Here we give the tree diagram of all strings starting with 1:



Now note that, the upper bound of Theorem 2 tells us we can give at most three 1-transpositions when n is even (i.e. $n-k$ is odd) and two 1-transpositions when n is odd (i.e. $n-k$ is even). Note that, if we could give a 2-transposition at each step, we would get the bound of $\lceil \frac{n-k}{2} \rceil$. For a n length string, if we can give a 2 transposition, the resulted reduced string may start with 1, 0 or 2. For the latter two cases, we can use relabeling as mentioned before. Therefore we can safely state that the reduced string will have any of the 8 prefixes of List 1.

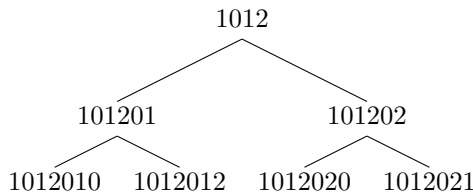
Hence, it suffices to prove the bound considering each of the prefixes of List 1. We will now follow the following strategy:

We will take each of the prefixes of List 1 and expand it (by adding symbols) to construct all possible strings of length greater than or equal to 7. Strictly speaking, we will not consider all possible strings; rather we will continue to expand until we get a 2-transposition, since afterwards, any further expansion would also guaranty a 2-transposition. Suppose s is one such string. We will take s and try to give a 2-transposition with any of its prefix. If we succeed, then, clearly, we are moving towards the best case and we only need to work with the reduced string. If we can't give a 2-transposition, we specifically deal with s and show that the bound holds. Now if we can give a 2-transposition, the reduced string will have any of the 8 prefixes (using relabeling if needed) and we will show that all strings of these cases will follow the bound.

Firstly it is easy to note that, the prefixes 1010 and 1212 themselves have 2-transpositions (Lemma 2). Therefore, we can safely exclude them from the following discussion. In what follows, when we refer to the prefixes of List 1, we would actually mean all the prefixes excluding 1010 and 1212. Now, to expand the prefixes, if we add 10 or 12, all of them would be able to give a 2-transposition (Lemma 2). Therefore, in what follows, we consider the other cases. Now we analyze each of the prefixes below.

1012

We first give the tree diagram of all string having prefix 1012



If we add 01, we can only add (By 'add' we mean append) 0 or 2. The resulting expanded string becomes 1012010 or 1012012. In both cases, we can give a 2-transposition. On the other hand, if we add 02, we can add 0 or 1 next and the string becomes 1012020 or 1012021. The string 1012020 satisfies the bound as follows: $\overline{1012020} \Rightarrow 012020 \Rightarrow \overline{012020} \Rightarrow 0120 \Rightarrow \overline{0120} \Rightarrow 120$. Here, $n = 7$ and we need only 3-transpositions holding the bound true. Now if we add 1, the string becomes 10120201. For this one as well the bound holds as follows: $\overline{10120201} \Rightarrow 0120201 \Rightarrow \overline{0120201} \Rightarrow 20201 \Rightarrow \overline{20201} \Rightarrow 201$. Here, $n = 8$ and we needed 3-transpositions. Now, it can be easily checked that strings like 1012020(20)* or 1012020(20)*1 the bound holds using the same strategy

as shown above. Adding anything with 1012020(20)*1 will also give a 2-transposition as follows. Clearly, we first need to add either 0 or 2 and immediately Lemma 2 would apply.

Now we consider 1012021. The bound holds for this one as well as follows: $\overline{1012021} \Rightarrow 012020 \Rightarrow \overline{012021} \Rightarrow 0121 \Rightarrow \overline{0121} \Rightarrow 201$. Now, strings like 1012(02)*1 can also be handled similarly hence the bound holds for them as well. Adding anything with 1012(02)*1 will also give a 2-transposition (Lemma 2).

1021

This prefix is ending with 1 and adding anything will give a 2-transposition (Lemma 2).

1020

We first add 21 with this prefix. Then, adding anything with 102021 will give a 2-transposition (Lemma 2). Now, for 102(02)⁺1, we need a 1-transposition to move the initial 1 at the end. Now the upper bound holds because the rest of the string is binary (Theorem 1). Also, adding anything with 102(02)⁺1 will give a 2-transposition (Lemma 2).

Now, if we add 20 with the prefix we get 102020. Next we add 1 or 2 and get 1020201 or 1020202 respectively. For 1020201, the bound holds as follows: $\overline{1020201} \Rightarrow 020201 \Rightarrow \overline{020201} \Rightarrow 0201 \Rightarrow \overline{0201} \Rightarrow 201$. All strings like 1020(20)⁺1 can also be handled similarly and hence the bound holds for them as well. And adding anything with 1020(20)⁺1 will give a 2-transposition (Lemma 2).

Now for 1020202, the bound holds as follows: $\overline{1020202} \Rightarrow 210202 \Rightarrow \overline{210202} \Rightarrow 2102 \Rightarrow \overline{2102} \Rightarrow 102$. Now, string like 10202(02)⁺ can be handled similarly and hence the bound holds. For 10202(02)⁺1, we need a 1-transposition to move the initial 1 at the end. Now the upper bound holds because the rest of the string is binary (Theorem 1). Also adding anything with 10202(02)⁺1 will give a 2-transposition (Lemma 2).

1201

This prefix is ending with a 1 and adding anything will give a 2-transposition (Lemma 2).

1202

Here, we can employ relabeling and map 2 to 0 and 0 to 2 to get 1020. Now recall that we have already considered 1020 before and hence we are done.

1210

Here we can again employ relabeling and map 2 to 0 and 0 to 2 to get 1012. And since we have already considered 1012, we are done.

This completes the proof. \square

4 Algorithm to group fully binary and ternary strings

Algorithm 1: GroupByPrefixTransposition
(s :input string)

Input: s , a fully binary or ternary string

initialization;
 $k \leftarrow 2$ if s is binary;
 $k \leftarrow 3$ if s is ternary;
 $count \leftarrow 0$;
 $twoTranspDone \leftarrow false$;

while $|s| > k$ **do**

for $i = 1; i < |s|; i = i + 1$ **do**

take the first symbol of input string ;
take the i^{th} symbol of the input string ;
append these two symbols ;
call this string $temp$;
check whether this string is a substring of
the current suffix ;

for $j = i + 1; j < (|s| - 1); j = j + 1$ **do**

take the substring named *consecutive*
of input string from j to $j + 2$;
if $consecutive == temp$ **then**

perform a 2-transposition ;
 $count \leftarrow count + 1$;
 $twoTranspDone \leftarrow true$;
break ;

end

end

if $twoTranspDone == false$ **then**

perform a 1-transposition using
Algorithm 2 ;
 $s \leftarrow$ return value from 2 ;
 $count \leftarrow count + 1$;

end

end

$twoTranspDone \leftarrow false$;

end

In this section we present Algorithm 1 to group fully binary and ternary strings which satisfies the proposed bound. We take all possible prefixes and try to find a 2-transposition in the suffix. If no 2-transposition is found, the algorithm gives a 1-transposition with the current prefix. It runs until the length is 2 and 3 for fully binary and fully ternary strings respectively. A variable count is initialized to zero and after the algorithm runs, it contains total number of prefix transpositions required. Algorithm 1 uses Algorithm 2 to perform 1 transposition (which is always available as discussed in Lemma 1). Algorithm 2 is simple, it first checks if the first character exists in the suffix. If so, it performs a 1-transposition. If not, it will increase the length of prefix and checks if the last character of the prefix exists in the suffix. If so, it performs a 1-transposition with that prefix and suffix.

Before formally proving the correctness of our algorithm, it will be useful to provide some classifications of the fully binary and ternary strings. In the following section, we will first classify the fully binary and ternary strings along two different dimensions and then discuss

the correctness of our algorithm. Notably, apart from being useful in proving the correctness of our algorithm, the classifications provided below would be interesting in their own right.

Algorithm 2: Do1Transposition (s :input string)

Input: s is a fully binary or ternary string passed as a parameter from Algorithm 1

$forward \leftarrow false$;
take the suffix of input string leaving the initial symbol of s ;
if the first symbol exists in the suffix **then**

perform a 1-transposition to s ;
 $forward \leftarrow true$;

end

if $forward == false$ **then**

for $p = 1; p < (|s| - 1); p = p + 1$ **do**

take a string suffix, substring of input
string from $(p + 1)$ to end ;
if p^{th} symbol occurs at suffix **then**

perform a 1-transposition ;

end

end

end

$return s$;

5 Classification

5.1 Classes of Ternary Strings

In this section, we will identify two classes of fully ternary strings. We have already shown that it needs at most $\lceil \frac{n-k}{2} \rceil + 1$ prefix transpositions to group ternary strings. We will discuss two classes; the first one consists of the strings which requires $\lceil \frac{n-k}{2} \rceil$ prefix transpositions (Class 1) and the other will need $\lceil \frac{n-k}{2} \rceil + 1$ (Class 2). We are not able to classify all the fully ternary strings; however below we provide classification of a significant number of strings:

- All length four strings satisfy $\lceil \frac{n-k}{2} \rceil$ bound, so we put them in Class 1.
- $(10)^+102$, $(10)^+210$, $(10)^+212$, $(12)^+010$, $(12)^+012$, $(12)^+120$ belong to Class 1. We can apply relabeling on these to find strings starting with 0 or 2 and they also belong to Class 1.
- If a ternary string is reduced to any one of the previous strings by a series of 2-transpositions, then that particular string also belongs to Class 1.

In Class 2, we put the strings which need $\lceil \frac{n-k}{2} \rceil + 1$ prefix transpositions to be grouped.

- 10120, 10121, 10201, 10202, 12021, 12020, 12101, 12102 and those we get applying relabeling on them are in Class 2. There is no 2-transposition available in these strings.
- $(10)^+120$, $(10)^+121$, $(10)^+201$, $(10)^+2(02)^+$, $(12)^+021$, $(12)^+0(20)^+$, $(12)^+101$, $(12)^+102$, $(10)^+1201$, $(10)^+2021$, $(12)^+0201$, $(12)^+1021$ and

those which we get by applying relabeling on these strings.

- Ternary strings reduced to any of these by a series of 2-transpositions will also be placed in Class 2.

Now, we will categorize all the fully ternary strings from a different perspective. We call a string “Good string” when Algorithm 1 gives optimal result. All the strings in Class 1 and Class 2 specified earlier are “Good strings”. There are strings for which Algorithm 1 will not give optimal results. We name those strings as “Hard strings”. $(02)^+10202$, $(10)^+21010$, $(21)^+02121$ are examples of “Hard strings”. To elaborate, let us take a fully ternary string 0210202 as an example. Algorithm 1 will give the following set of prefix transpositions: $0\bar{2}10202 \Rightarrow \bar{1}0\bar{2}02 \Rightarrow \bar{2}102 \Rightarrow 102$. Here we need 3 prefix transpositions which satisfy the bound. However, an optimal sequence will need 2 prefix transpositions to group this string $0\bar{2}10202 \Rightarrow 0\bar{2}102 \Rightarrow 102$. So, 0210202 is a “Hard string”.

5.2 Correctness of Algorithm 1

With the above classification of the ternary strings at our hand, we are now ready to prove the correctness of our algorithm. We prove the following theorem:

Theorem 3. *Given a fully binary and ternary string, Algorithm 1 is always able to group satisfying the bound given in Theorems 1 and 2.*

Proof. The proof is simple for fully binary strings because Algorithm 1 will always give optimal results for fully binary strings. So we need only prove that Algorithm 1 will always satisfy the proposed bounds for ternary strings. Since, by definition, Algorithm 1 provides optimal results for “Good strings”, we only need to consider “Hard strings” and prove that Algorithm 1 will satisfy the upper bound while dealing with the “Hard strings”. Consider a hard string. Algorithm 1 may choose a 2-transposition which may not be optimal. Then it will reduce the string into any one of the followings: $(10)^+120$, $(10)^+121$, $(10)^+201$, $(10)^+2(02)^+$, $(12)^+021$, $(12)^+0(20)^+$, $(12)^+101$, $(12)^+102$, $(10)^+1201$, $(10)^+2021$, $(12)^+0201$, $(12)^+1021$ or any fully ternary string of length 4 and those which we get by applying relabeling any of these. These strings satisfy the upper bound. So, overall the Algorithm 1 will satisfy the upper bound. \square

6 Sorting

The sorting distance $d_s(s)$ of a fully k ary string s is defined as the minimum number of prefix transposition required to sort the string to one of length k . We again consider normalized strings.

6.1 Sorting Binary Strings

Theorem 4. (Bound for Binary strings) *Let s be a fully binary string. Then, $d_s(s) \leq \lceil \frac{n-k}{2} \rceil$.*

Proof. As binary strings have only 2 letters, after grouping they are already sorted (in ascending or descending order). So, the upper bound is $d_s(s) \leq \lceil \frac{n-k}{2} \rceil$. \square

6.2 Sorting Ternary Strings

Theorem 5. (Bound for Ternary strings) *Let s be a fully Ternary string. Then, upper bound for sorting ternary string is $d_g(s) \leq \lceil \frac{n-k}{2} \rceil + 2$.*

Proof. After grouping a ternary string, we have the following grouped strings: 012, 021, 102, 120, 201 and 210. Among these, 012 and 210 are already sorted. We need one more 0 transposition to sort 021, 102, 120 and 201. Hence the result follows. \square

6.3 Algorithm to sort fully binary and ternary strings

We first apply Algorithm 1 of Section 4 on the input string. Then, if the string is binary, it is already sorted. If it is ternary and the after grouping we have any one of 021, 102, 120 and 201, we need 1 more prefix transposition to sort the string.

Algorithm 3: SortByPrefixTransposition (s :input string)

```

Input:  $s$  is a fully binary or ternary string
;
run Algorithm 1 on  $s$  ;
if  $s$  is binary then
     $s$  is sorted after grouping ;
    finish ;
end
else
    if  $s$  is in 021, 102, 120, 201 then
        perform one 0-transposition to sort  $s$  ;
    end
    else
         $s$  is already sorted ;
    end
end

```

7 Conclusion

In this paper, we have discussed grouping and sorting of fully binary and ternary strings when the allowed operation is prefix transposition. Following the work of [4], we have handled grouping by prefix transpositions of binary and ternary strings first and extended the results for sorting. In particular we have proved that, for binary strings the grouping distance is $d_g(s) = \lceil \frac{n-k}{2} \rceil$ and for ternary string we have $\lceil \frac{n-k}{2} \rceil \leq d_g(s) \leq \lceil \frac{n-k}{2} \rceil + 1$, where n is the length of the string and k is the arity. On the other hand, for sorting binary and ternary strings the sorting distance $d_s(s)$ is upper bounded by $\lceil \frac{n-k}{2} \rceil$ and $\lceil \frac{n-k}{2} \rceil + 2$ respectively. As has already been mentioned, we are now considering the higher-arity alphabets as an extension of our work.

References

- [1] X. Chen, J. Zheng, Z. Fu, P. Nan, Y. Zhong, S. Lonardi, and T. Jiang. Assignment of orthologous genes via genome rearrangement. *IEEE/ACM Trans. Comput. Biology Bioinform.*, 2(4):302–315, 2005.
- [2] D.A. Christie and R.W. Irving. Sorting strings by reversals and by transpositions. *SIAM J. Discrete Math.*, 14(2):193–206, 2001.
- [3] H. Eriksson, K. Eriksson, J. Karlander, L.J. Svensson, and J. Wästlund. Sorting a bridge hand. *Discrete Mathematics*, 241(1-3):289–300, 2001.
- [4] C.A.J. Hurkens, L. van Iersel, J. Keijsper, S. Kelk, L. Stougie, and J. Tromp. Prefix reversals on binary and ternary strings. *SIAM J. Discrete Math.*, 21(3):592–611, 2007.
- [5] A.J. Radcliffe, A.D. Scott and E.L. Wilmer. Reversals and transpositions over finite alphabets. *SIAM J. Discrete Math.*, 2006.