

# Block Sorting is Hard

Wolfgang W. Bein \*

Lawrence L. Larmore †

Shahram Latifi ‡

I. Hal Sudborough §

## Abstract

*Block sorting is used in connection with Optical Character Recognition (OCR). Recent work has focused on finding good strategies which work in practice.*

*In this paper, we show that optimizing block sorting is  $\mathcal{NP}$ -hard. Along with this result, we give new non-trivial lower bounds. These bound can be computed efficiently. We define the concept of “Local Property Algorithms” and show that several previously published block sorting algorithms fall into this class.*

## 1 Preliminaries

Consider an ordered alphabet  $\Sigma$  with ordering relation “ $<$ ”. For a fixed string  $w$  over  $\Sigma$ , if symbol  $a$  has position  $i$  in  $w$  and symbol  $b$  has position  $j$ , for  $i < j$ , we write  $a \dots b$ . (Thus the relation “ $\dots$ ” is an order relation which depends on  $w$ .) We call “ $\dots$ ” the *position order*. Two symbols  $a, b$  in  $w$  are called *consecutive* if  $a$  has position  $i$  in  $w$  and  $b$  has position  $i + 1$ , and  $a < b$ , and there is no  $c$  in  $w$  with  $a < c < b$ . A string  $w$  is called *irreducible* if there are no consecutive elements  $a, b$  in  $w$ .

In block sorting we are only interested in irreducible strings. Starting with an irreducible string we are interested in the minimum number of moves to sort the string: A *move* is defined as deleting a symbol of  $w$  at position  $i$  while inserting it at position

$j \neq i$ . We write  $move(x)$  for the move which deletes the symbol  $x$ . If the resulting string has consecutive symbols, they are joined into one block, which is then considered one symbol. For convenience, after such a move we always name the block with a symbol from the block not involved in the move to symbolize the joined block. The block sorting problem is, given a string  $w$ , what is the minimum number of moves to sort the string.

Block sorting is motivated by applications in Optical Character Recognition; see [6, 10, 11]. Text regions, referred to as *zones* are selected by drawing rectangles around them. Here the order of zones is significant, but in practice the output generated by a zoning procedure may be different from the correct text. Moving the pieces to the correct order corresponds to a block sorting problem. To evaluate the performance of a given zoning procedure it is of interest to find the minimum number of moves needed to obtain the correct string from the zones generated by the zoning procedure. Recently, Latifi *et al.* [6] have performed various experiments to test a number of strategies that seem to perform well in practice.

Unfortunately, as we shall see in this paper, block sorting is inherently hard. In fact, no better than a 3-competitive approximation algorithm is known. Sorting problems under various operations have been studied extensively. We mention work on sorting with prefix reversals [3, 5, 9], transpositions [1, 2] and block moves [7, 8]. In Section 2 we give a number of non-trivial lower bounds. In Section 3, we show that optimizing block sorting is  $\mathcal{NP}$ -hard. In Section 4 we examine a number of strategies that have been studied previously in connection with OCR [6], and show that these strategies can all be viewed as what we call “Local Property Algorithms” and that the run times of their implementations can be substantially improved.

\*Department of Computer Science, University of Nevada, Las Vegas, NV 89154-4019. Email: [bein@cs.unlv.edu](mailto:bein@cs.unlv.edu). Research supported by NSF grant CCR-9821009.

†Department of Computer Science, University of Nevada, Las Vegas, NV 89154-4019. Email: [larmore@cs.unlv.edu](mailto:larmore@cs.unlv.edu). Research supported by NSF grant CCR-9821009.

‡Department of Electrical and Computer Engineering, University of Nevada, Las Vegas, NV 89154-4026. Email: [latifi@ee.unlv.edu](mailto:latifi@ee.unlv.edu).

§Department of Computer Science, University of Texas at Dallas, Richardson, TX 75083-0688 Email: [hal@utdallas.edu](mailto:hal@utdallas.edu).

## 2 Lower Bounds

Given a string  $w$  let  $\hat{w}$  be the string consisting of the same symbols, sorted under the order “ $<$ ”. The dual block sorting problem is: What is the minimum number of moves to sort the string  $\hat{w}$  under position ordering “ $\dots$ ” (of  $w$ .) We refer to  $\hat{w}$  under position ordering as the *dual string*. Throughout this section we make use of a number of straightforward lemmas, whose proofs will be given in the full paper.

**Lemma 1** *The minimum number of moves for block sorting  $w$  under order “ $<$ ” is equal to the minimum number of moves to block sort  $\hat{w}$  under order “ $\dots$ ” of  $w$ .*

We say there is a *blue edge* from  $a$  to  $b$  if  $a > b$  and  $b$  is immediately after  $a$  in  $w$ , i.e.,  $ab$  is a substring of  $w$ . We also refer to a blue edge as a *reversal*. Alternatively, we say there is a *green edge* if  $a \dots b$  and  $a > b$  and there is no  $c$  in  $w$  with  $a > c > b$ . A green edge is also called a *dual reversal*. Note that a green edge implies a reversal in the dual string, and vice versa.

**Lemma 2** *Let  $w$  be a string with  $k$  blue edges and  $\ell$  green edges. After the removal of any symbol the number of blue edges  $k'$  and green edges  $\ell'$  of the resulting string  $w'$  does not increase, i.e.  $k' \leq k$  and  $\ell' \leq \ell$ . Furthermore, if  $|w'| = n - 2$ , then  $k' = k - 1$  or  $\ell' = \ell - 1$ , and if  $|w'| = n - 3$ , then  $k' = k - 1$  and  $\ell' = \ell - 1$ .*

**Lemma 3** *Let  $w$  be a string of length  $n$  with  $k$  blue edges and  $\ell$  green edges. Then  $m \geq n - k - \ell - 1$ , where  $m$  is the minimum number of moves needed to block sort  $w$ .*

Combining lemma 2 and 3, we obtain the following lower bound.

**Theorem 1** *Given a string  $w$  of length  $n$  with  $k$  reversals and  $\ell$  dual reversals, Then the minimum number of moves  $m$  required to block sort  $w$  is bounded from below by*

$$m \geq \max \begin{cases} n - k - \ell - 1 \\ k \\ \ell \end{cases}$$

By averaging the three components in the inequality of theorem 1, we obtain:

**Corollary 1** *For a string  $w$  of length  $n$  the the minimum number of moves  $m$  required to block sort  $w$  is bounded from below by  $m \geq \lceil \frac{n-1}{3} \rceil$ .*

We say that a sorting is *perfect* if the number of moves is equal to the number of blue edges. Note that Theorem 1 implies that a perfect sorting is optimal. Though the converse does not hold, it is useful to characterize perfect sortings. To this end, fix a specific block sorting of  $w$ .

We say that there is a *red edge* from  $a$  to  $b$  if

- $a < b$  and  $a \dots b$
- $a$  and  $b$  are joined before either is moved.
- If  $a \dots c \dots b$ , then  $c$  is moved before  $a$  and  $b$  are joined.

The *red-blue graph* is the graph whose nodes are the symbols in  $w$  and whose edges are the red edges and the blue edges.

The following lemmas show that red edges are essentially “savings.” As before, we omit the proofs here.

**Lemma 4** *The number of moves plus the number of red edges is one less than the length of  $w$ .*

**Lemma 5** *The red-blue graph is acyclic.*

From Lemmas 4 and 5 we immediately obtain:

**Theorem 2** *The sorting is perfect if and only if the red-blue graph is a tree.*

We now give a number of results (again without proof) which limit the number of red edges in the red-blue graph and which are useful later.

**Lemma 6** *For any given node  $x$ , there can be at most one  $y < x$  for which there is a red edge from  $y$  to  $x$ . Furthermore, there can be at most one  $z > x$  for which there is a red edge from  $x$  to  $z$ .*

**Lemma 7** *If  $a \dots c \dots b \dots d$ , then there cannot be both a red edge from  $a$  to  $b$  and a red edge from  $c$  to  $d$ .*

The following lemma follows directly from duality:

**Lemma 8** *If  $a < c < b < d$ , then there cannot be both a red edge from  $a$  to  $b$  and a red edge from  $c$  to  $d$ .*

**Lemma 9** *If  $a < c < d < b$  and  $c\dots a\dots b\dots d$ , then there cannot be both a red edge from  $a$  to  $b$  and a red edge from  $c$  to  $d$ .*

We conclude this section with an alternative lower bound, which is based on counting alternating cycles in a certain directed graph. We call  $w$  *anchored* if the first symbol of  $w$  is also the smallest, and the last symbol is the largest. If not, we can simply “anchor”  $w$  by adding the prefix 0, defined to be smaller than any symbol of  $w$ , and appending “ $\infty$ ,” defined to be smaller than any symbol of  $w$ . Trivially, the number of steps required for block sorting does not change. Henceforth in this section, we assume that  $w$  is anchored. As before, let  $n$  be the length of  $w$ .

We now define *black* and *white* edges. If  $ab$  is a substring of  $w$ , we say that there is a white edge from  $b$  to  $a$ . If  $a < b$  and there is no symbol of  $w$  which is larger than  $a$  and smaller than  $b$ , we say that there is a black edge from  $a$  to  $b$ . We define the *breakpoint* graph  $G$  to be the directed graph whose nodes are the symbols of  $w$  and whose edges are the black and white edges. We define an *alternating cycle* (or just “*cycle*” if the context is clear) to be a directed cycle whose edges are alternately white and black. Note that every symbol except the first and last has black out-degree 1, white out-degree 1, black in-degree 1, and white in-degree 1. The first symbol has black out-degree 1 and white in-degree 1, while the last symbol has black in-degree 1 and white out-degree 1. Thus, each edge of  $G$  lies in precisely one alternating cycle.

We define the *black length* of a cycle to be its number of black edges. We say that a cycle is *odd* if it has odd black length, otherwise, it is *even*. Let  $odd\_cycles(w)$  be the number of odd cycles in  $G$ . We note that the black length of any cycle is at least 2, since we assume there are no consecutive symbols in  $w$ .

**Theorem 3** *Let  $w$  be an anchored string of length  $n$ . The number of moves required to sort  $w$  is at least  $\frac{1}{2}(n - 1 - odd\_cycles(w))$ .*

*Proof:* If  $C$  is any alternating cycle of the breakpoint graph  $G$  of  $w$ , and the black length of  $C$  is  $c$ , we define  $\Phi(C) = \lfloor \frac{c}{2} \rfloor$ . Thus,  $\Phi(C) = 1$  if  $c = 3$ , while  $\Phi(C) = 2$  if  $c = 4$  or  $c = 5$ , and so forth. We define  $\Phi(w) = \sum_C \Phi(C)$ , where  $C$  ranges over all alternating cycles of  $G$ . Note that  $\Phi(w) = \frac{1}{2}(n - 1 - odd\_cycles(w))$ .

We can show that if  $x$  is any symbol other than the first or last symbol, and if  $w'$  is the resulting string after  $move(x)$ , then  $\Phi(w) - \Phi(w') \leq 1$ . Therefore, the number of moves needed to sort  $w$  is at least  $\Phi(w)$ .

Cycles are edge disjoint, but not node disjoint. We define a Type I symbol to be a symbol which appears as a node twice in one cycle, and a Type II symbol to be a symbol which appears in two different cycles. Each symbol, other than the first and last symbols, is either Type I or Type II. Recall that the move  $move(x)$  is defined by choosing a symbol  $x$  of  $w$ , deleting  $x$ , then combining all consecutive symbols, if any.

We summarize the effects of  $move(x)$  on the numbers of even and odd cycles, as follows:

1. Suppose  $x$  is a Type II symbol. Let  $P$  and  $Q$  be the two cycles that contain  $x$ . Suppose  $P$  has black length  $p$  and  $Q$  has black length  $q$ . The effect of  $move(x)$  on the alternating cycles is: (a)  $P$  and  $Q$  are removed. (b) There is a new cycle  $T$  whose black length is  $p + q - 1$ . and (c) All other cycles remain.
2. Suppose  $x$  is a Type I symbol, which appears twice in a cycle  $T$ , of black length  $t$ . (Then the actual length of the  $T$  is  $2t$ .) Then  $T$  is the union of two cycles (in this case, not alternating cycles) which meet at  $x$ , of lengths  $u$  and  $v$ , where  $u$  and  $v$  are both odd and  $u + v = 2t$ . Note, also, that  $u, v \geq 3$ . The effect of  $move(x)$  on the alternating cycles is as follows: (a)  $T$  is removed. (b) If  $u \geq 5$ , then there is a new cycle  $U$  of black length  $\frac{u-1}{2}$ . (c) If  $v \geq 5$ , then there is a new cycle  $V$  of black length  $\frac{v-1}{2}$ . (d) All other cycles remain.

Checking all the cases, we can verify that the value of  $\Phi$  cannot decrease by more than 1, and we are done.  $\square$

### 3 The Block Sorting Problem is $\mathcal{NP}$ -Complete

The 0-1 block sorting problem is, given string  $w$ , can  $w$  be block sorted in  $k$  moves? Henceforth, when we say the “block sorting problem,” or just “block sorting,” we mean the 0-1 block sorting problem. We show now that this problem is  $\mathcal{NP}$ -complete.

Trivially, block sorting is in  $\mathcal{NP}$ . To prove  $\mathcal{NP}$ -completeness, we reduce *boolean 3-satisfiability*, a well-known [4]  $\mathcal{NP}$ -complete problem, to block sorting.

Suppose  $\mathcal{B} = \mathcal{C}^1 \mathcal{C}^2 \dots \mathcal{C}^m$  is a boolean expression consisting of the conjunction of  $m > 0$  clauses, each of which is the disjunction of exactly 3 terms, each of which is either a variable or the negation of a variable. Assume there are  $n$  variables. We introduce an ordered alphabet  $\Sigma_{n,m}$  consisting of  $4nm + 2m + 4n + 1$  distinct symbols. We then construct a string  $\mathcal{X}$  of length  $8m + 4n + 1$  over  $\Sigma_{n,m}$  in which no symbol is repeated. Finally, we show that  $\mathcal{B}$  is satisfiable if and only if  $\mathcal{X}$  can be block sorted in  $6m + 2n - 1$  moves.

$\Sigma_{n,m}$  consists of the following symbols:

- $p_i^j, \bar{p}_i^j, q_i^j, \bar{q}_i^j$ , for all  $1 \leq i \leq n$  and  $1 \leq j \leq m$ . We call these *term symbols*. We call  $p_i^j$  and  $\bar{p}_i^j$  *left term symbols*, and we call  $q_i^j$  and  $\bar{q}_i^j$  *right term symbols*.
- $\ell^j, r^j$  for all  $1 \leq j \leq m + n$ . We call these *clause control symbols*
- $u_i, v_i$  for all  $1 \leq i \leq n$ . We call these *variable control symbols*
- $s$ . We call this the *separator*.

The ordering on  $\Sigma_{n,m}$  is generated by the following rules:

- $u_i < p_i^k < \bar{p}_i^k < q_i^j < \bar{q}_i^j < v_i < s$  for all  $1 \leq i \leq n$  and  $1 \leq j < k \leq m$ .
- $v_{i-1} < u_i$  for all  $1 < i \leq n$
- $s < \ell^k < r^k < \ell^j < r^j$  for all  $1 \leq j < k \leq m + n$ .

Assume now that for an instance of the boolean 3-satisfiability problem the names of the variables are  $x_i$  for  $1 \leq i \leq n$ . A clause must be of the form  $z_a + z_b + z_c$ , where each  $z_i$  is either  $x_i$  or  $\bar{x}_i$ . Without loss of generality,  $a > b > c$ . We first define the *simple encoding* of a clause, using symbols which have no superscripts. The simple encoding of  $x_i$  uses the two symbols  $p_i$  and  $q_i$ , while the simple encoding of  $\bar{x}_i$  uses the two symbols  $\bar{p}_i$  and  $\bar{q}_i$ . The simple encoding of a clause is eight symbols, beginning with  $\ell$  and ending with  $r$ . The remaining six symbols are those needed for the terms, but interleaved; first the  $p$ 's, then the  $q$ 's. For example, the simple encoding of the clause  $(x_5 + x_3 + \bar{x}_2)$  is  $\ell p_5 p_3 \bar{p}_2 q_5 q_3 \bar{q}_2 r$ .

For the true encoding of a clause, insert the index of that clause as a superscript for every symbol. For example, if  $\mathcal{C}^4 = (x_5 + x_3 + \bar{x}_2)$ , then its encoding is  $\ell^4 p_5^4 p_3^4 \bar{p}_2^4 q_5^4 q_3^4 \bar{q}_2^4 r^4$ .

After the true encoding of the clauses we add "control sequences"  $\ell^{m+i} u_i v_i r^{m+i}$ , for  $1 \leq i \leq n$ ,  $\ell^{m+i} u_i v_i r^{m+i}$ . Finally, we finish off the encoding:  $\mathcal{X}$  consists of  $s$ , followed by the encodings of the clauses in order, followed by the control sequences for each  $i$  from 1 to  $n$ .

We are now ready to prove the result: We first prove that satisfiability of a clause implies a perfect sorting.

**Lemma 10** *If  $\mathcal{B}$  is satisfiable, there exists a perfect sorting of  $\mathcal{X}$ .*

*Proof:* Given  $\mathcal{B}$ , and given a satisfying assignment of the variables, we now show how to block sort  $\mathcal{X}$  in exactly  $6m + 2n - 1$  steps.

For each  $1 \leq j \leq m$ , designate one term of  $\mathcal{C}^j$  which is true under the assignment. Call this the *relevant term* of  $\mathcal{C}^j$ . The steps are as follows:

1. Delete all symbols which form the encodings of terms which are not relevant. This takes  $4m$  steps.
2. For all  $j$ , starting with 1 and ending with  $m$ , in that order, join and the delete the two symbols which form the encoding of the relevant term of  $\mathcal{C}^j$ . These symbols will be either  $p_i^j q_i^j$  or  $\bar{p}_i^j \bar{q}_i^j$  for some  $i$ . This takes  $m$  steps. Note that elements become available for joining due to the order in which they are processed.
3. Join all pairs of the form  $u_i v_i$  and delete them. This takes  $n$  steps.
4. For all  $2 \leq j \leq m + i$ , join the symbols  $\ell^j r^j$  and delete them. This takes  $m + n - 1$  steps.

The remaining symbols are sorted.  $\square$

We now show that the existence of a perfect sorting implies satisfiability.

**Lemma 11** *If there is a perfect sorting of  $\mathcal{X}$ , then  $\mathcal{B}$  is satisfiable.*

*Proof:* If there is a perfect sorting, then by Theorem 2 the red-blue graph is connected, and acyclic, i.e. a tree. It has exactly  $6m + 2n - 1$  blue edges.

All blue edges connect symbols which are adjacent in  $\mathcal{X}$ . The components of the blue graph, which we call blue components, are thus substrings of  $\mathcal{X}$ . There are  $2m + 2n + 2$  blue components. We refer to the space between two blue components as a *gap*. The  $2m+2n+1$  red edges must connect those components, and must not violate the conditions given by Lemmas 8, 7, and 9. We shall show that these requirements imply the needed properties.

Claim A: There is no red edge from  $s$  to  $r^{\hat{j}}$ . Instead, for some  $1 \leq \hat{j} \leq m$ , there is a red edge from  $s$  to  $\ell^{\hat{j}}$ . Proof: If there is a red edge from  $s$  to  $r^{\hat{j}}$ , there can be no other red edge across the gap to the left of  $r^{\hat{j}}$ , by Lemmas 6 and 8. Thus, the graph is disconnected, a contradiction. On the other hand, there must be a red edge from  $s$  to some  $x$ , else the graph is disconnected. Since  $s < x$ ,  $x$  must be either  $\ell^j$  or  $r^j$ , for some  $j$ , by Lemma 7. Since there is no red edge from  $s$  to  $r^j$ , the claim follows.

Claim B: For all  $j$ , there is a red edge from  $\ell^j$  to  $r^j$ . Proof: First consider  $j > \hat{j}$ : If  $\hat{j} = n + m$ , there is nothing to prove. Otherwise, we proceed by induction on  $j$ , beginning with  $j = n + m$ . Since  $r^j$  is the first symbol of its component (*i.e.*, there is no connection via a blue edge) there must be a red edge from  $x$  to  $y$  where  $x \dots r^j$  and either  $y = r^j$  or  $r^j \dots y$ . The case  $r^j \dots y$  is ruled out by the inductive hypothesis, using Lemma 8. Thus there is a red edge from  $x$  to  $r^j$ . By Claim A and Lemma 7,  $s < x < r^j$ . The only choice for  $x$  is  $\ell^j$ .

Now we consider the case  $j < \hat{j}$ : For any  $j < \hat{j}$ , there is a red edge from  $\ell^j$  to  $r^j$ . If  $\hat{j} = 1$ , there is nothing to prove. Otherwise, we use induction on  $j$ , beginning with  $j = \hat{j} - 1$ . Suppose  $j = \hat{j} - 1$ . There must be a path from  $\ell^j$  to  $r^j$ . By Lemmas 7 and 8, and by Claim A, the only way such a path can exist is for there to be a red edge from  $\ell^j$  to  $r^j$ . Now suppose  $j < \hat{j} - 1$ . There must be a path from  $\ell^j$  to  $r^j$ . By the inductive hypothesis, there are red edges from  $\ell^k$  to  $r^k$  for all  $j < k < \hat{j}$ . By Lemmas 8 and 9, the only way such a path can exist is if there is a red edge from  $\ell^j$  to  $r^j$ .

For the case  $j = \hat{j}$ , we need to show the following claim first:

Claim C: There is a red edge from  $u_i$  to  $v_i$ . Proof: There must be a red edge connecting  $v_i$  to some other symbol. By the previous arguments in Claim B and Lemma 7, the only possibility for that other symbol is  $u_i$ . This concludes the proof of Claim C.

Returning to Claim B, we now can establish that

there is a red edge from  $\ell^{\hat{j}}$  to  $r^{\hat{j}}$ : Since  $r^{\hat{j}}$  is the first symbol of its component, there must be a red edge from  $x$  to  $y$  where  $x \dots r^{\hat{j}}$  and either  $y = r^{\hat{j}}$  or  $r^{\hat{j}} \dots y$ . The case  $r^{\hat{j}} \dots y$  is ruled out by Claim A and Lemma 8. Claim A rules out the possibility that  $x \dots \ell^{\hat{j}}$ . By Claim C and Lemma 7,  $x$  cannot be a term symbol. Thus,  $x = \ell^{\hat{j}}$ , the only remaining possibility.

Claim D: For each  $1 \leq j \leq m$ , there is some  $i$  for which there is a red edge from  $p_i^j$  to  $q_i^j$ , in which case  $x_i$  is a term of  $\mathcal{C}^j$ , or from  $\bar{p}_i^j$  to  $\bar{q}_i^j$ , in which case  $\bar{x}_i$  is a term of  $\mathcal{C}^j$ . Furthermore, if there is a red edge from  $p_i^j$  to  $q_i^j$ , then there is no red edge from  $\bar{p}_i^k$  to  $\bar{q}_i^k$ . Proof: We first show that for each  $1 \leq j \leq m$ , there is a red edge from  $x$  to  $y$  where  $\ell^j \dots x \dots y \dots r^j$ , and where  $x$  is a left term symbols and  $y$  is a right term symbol. This is because the blue component containing the three right term symbols between  $\ell^j$  and  $r^j$  must have a red edge. All other possibilities for  $x$  are ruled out by Claim B.

To prove that  $\mathcal{B}$  is satisfiable we select a validating term for each clause such that no two validating terms are contradictory. Let the validating term of clause  $\mathcal{C}^j$  be  $x_i$  if the red edge connects  $p_i^j$  and  $q_i^j$ , and  $\bar{x}_i$  if the red edge connects  $\bar{p}_i^j$  to  $\bar{q}_i^j$ . Claim D guarantees no that two of the validating terms contradict each other.  $\square$

Lemma 10 and 11 now establish our result:

**Theorem 4** *Block sorting is  $\mathcal{NP}$ -complete.*

## 4 Local Property Algorithms

A local property algorithm attempts to minimize the number of moves to block sort a string by deciding locally whether to delete an element. For example, an algorithm that seeks to eliminate blue edges must remember whether any given symbol is the left end point of a blue edge. A local property algorithm could look for other properties, but we require that the algorithm maintain the information at a symbol in constant time per node update.

Specifically, a local property algorithm maintains a data structure called a *quad list*. If the string to be sorted is  $w = w_1, \dots, w_n$  then the quad list is a doubly linked list consisting of the elements  $(w_1, 1), \dots, (w_n, n)$  with a head pointer pointing to  $(w_1, 1)$ . Additionally there are bidirectional pointers for the dual order, *i.e.* there is a bidirectional pointer between  $(u, i)$  and  $(v, j)$  if  $u < v$  and there

exists no  $y$  in  $w$  such that  $u < y < v$ . We call these additional pointers the *value pointers*, the former *positional pointers*. Each node therefore stores four pointers. We call the four associated pointer fields, `left`, `right` for the positional pointers and `up`, `down` for value pointers. Thus for node  $x$  the previous or next elements in respective orderings are `x.left`, `x.right`, `x.down`, `x.up`, where by a slight abuse of notation we use the pointer to mean also the node under the pointer. Finally, we assume that we have a constant number of extra boolean fields to store local information at each node of the list as needed. For example, there could be an extra field `x.blue_edge` whose value is true if the symbol in  $x$  is the left end point of a blue edge. We require that such properties are local in the sense that the field can be calculated and maintained in  $O(1)$  time per node. The quad list can be set up in  $O(n \log n)$  time. Thus, a local property algorithm take  $O(n \log n)$  setup-time and  $O(n)$  for all subsequent steps.

Table 1 summarizes various useful local properties. A node  $x$  has property A if  $x$  is simultaneously the left end point of a blue edge and a green edge, or equivalently `x.right.down = x`. This property is useful, since any optimal algorithm can without loss of optimality delete such elements first. Lemma 2 suggests looking for nodes the deletion of which causes the string length to decrease by up to three. Properties B and C identify such nodes. If there are no nodes of type B and C then a possible strategy is to identify nodes the deletion of which will cause situations B and C in the next step; properties D, E, and F identify such nodes.

Type	Condition	Effect of Deletion
A	<code>x.r.d = x</code>	$n' \leq n - 1$
B	<code>x.d.r.d = x</code>	$n' \leq n - 2$
C	<code>x.l.u.l = x</code>	$n' \leq n - 2$
D	<code>x.l.u.l.l = x</code>	$n' \leq n - 1$ , <code>x.r</code> becomes type C
E	<code>x.d.r.d.d = x</code>	$n' \leq n - 1$ , <code>x.u</code> becomes typeB
F	<code>x.l.u.u.l = x</code>	$n' \leq n - 1$ <code>x.l.u</code> becomes type B

Table 1: Properties of node  $x$ , and the effect of deleting  $x$ . Note that `right`, `left`, `up`, `down` are abbreviated by `r`, `l`, `u`, and `d`.

Latifi *et al.* [6] give an  $O(n^2)$  algorithm, called

the 1-lookahead algorithm, for which they show that computer experimentation gives relatively good results in practice. Their algorithm is the local property algorithm with the properties of Table 1, where properties A ... F have alphabetical priority. If no property is available at any step the 1-lookahead algorithm simply arbitrarily removes the last node. Thus our data structures improves their run time to  $O(n \log n)$ .

## References

- [1] V. Bafna and P.A. Pevzner. Sorting by transposition. *SIAM Journal on Discrete Mathematics*, 11:224–240, 1998.
- [2] V. Bafna and P.A. Pevzner. Genome rearrangements and sorting by reversals. *SIAM Journal on Computing*, 25:272–289, 1999.
- [3] A. Caprara. Sorting by reversals is difficult. In *Proceedings 1st Conference on Computational Molecular Biology*, pages 75–83. ACM, 1997.
- [4] M. R. Garey and D. S. Johnson. *Computers and intractability – A guide to the theory of NP-completeness*. Freeman, San Francisco, 1979.
- [5] W. H. Gates and C. H. Papadimitriou. Bounds for sorting by prefix reversal. *Discrete Mathematics*, 27:47–57, 1979.
- [6] R. Gobi, S. Latifi, and W. W. Bein. Adaptive sorting algorithms for evaluation of automatic zoning employed in OCR devices. In *Proceedings of the 2000 International Conference on Imaging Science, Systems, and Technology*, pages 253–259. CSREA Press, 2000.
- [7] L.S. Heath and J.P.C. Vergara. Sorting by bounded block-moves. *Discrete Applied Mathematics*, 88:181–206, 1998.
- [8] L.S. Heath and J.P.C. Vergara. Sorting by short block-moves. *Algorithmics*, 28 (3):323–354, 2000.
- [9] H. Heydari and I. H. Sudborough. On the diameter of the pancake network. *Journal of Algorithms*, 25(1):67–94, 1997.
- [10] J. Kanai, S. V. Rice, and T.A. Nartker. Automatic evaluation of OCR zoning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 17 (1):86–90, 1995.
- [11] S. Latifi. How can permutations be used in the evaluation of automatic zoning evaluation? In *ICEE 1993*. Amir Kabir University, 1993.