

## Sorting by reversals and the theory of 4-regular graphs

Non Peer-reviewed author version

BRIJDER, Robert (2017) Sorting by reversals and the theory of 4-regular graphs. In: Theoretical computer science, 701, p. 40-53.

DOI: 10.1016/j.tcs.2017.02.033

Handle: <http://hdl.handle.net/1942/25631>

# Sorting by Reversals and the Theory of 4-Regular Graphs

Robert Brijder<sup>1</sup>

*Hasselt University, Belgium*

---

## Abstract

We show that the theory of sorting by reversals fits into the well-established theory of circuit partitions of 4-regular multigraphs (which also involves the combinatorial structures of circle graphs and delta-matroids). In this way, we expose strong connections between the two theories that have not been fully appreciated before. We also discuss a generalization of sorting by reversals involving the double-cut-and-join (DCJ) operation. Finally, we also show that the theory of sorting by reversals is closely related to that of gene assembly in ciliates.

*Keywords:* sorting by reversals, sorting by DCJ operations, genome rearrangements, 4-regular graphs, local complementation, gene assembly in ciliates

---

## 1. Introduction

Edit distance measures for genomes can be used to approximate evolutionary distance between their corresponding species. A number of genome transformations have been used to define edit distance measures. In this paper we consider the well-studied chromosome transformation called *reversal*, which is an inversion of part of a chromosome [34]. If two given chromosomes can be transformed into each other through reversals, then the difference between these two chromosomes can be represented by a permutation, where the identity permutation corresponds to equality. As a result, transforming one chromosome into the other using reversals is called *sorting* by reversals. The *reversal distance* is the *least* number of reversals needed to accomplish this transformation (i.e., to sort the permutation by reversals). In [27] a formula is given for the reversal distance, leading to an efficient algorithm to compute reversal distance. The proof of that formula uses a notion called the breakpoint graph. Subsequent streamlining of this proof led to the introduction of additional notions [31] such as the overlap graph and a corresponding graph operation. A reversal can be seen as

---

*Email address:* robert.brijder@uhasselt.be (Robert Brijder)

<sup>1</sup>Postdoctoral fellow of the Research Foundation – Flanders (FWO).

a special case of a double-cut-and-join (DCJ) operation that can operate either within a chromosome or between chromosomes. Similar as for reversals, one can define a notion of DCJ distance, and a formula for DCJ distance is given in [4].

The theory of circuit partitions of 4-regular multigraphs was initiated in [33] and is currently well-developed with extensions and generalizations naturally leading into the domains of, e.g., linear algebra and matroid theory. In this paper we show that this theory can be used to study sorting by reversals, and more generally sorting by DCJ operations. Moreover, we show how various aspects of the theory of circuit partitions of 4-regular multigraphs relate to the topic of sorting by reversals. In particular, we show that the notion of an overlap graph and its corresponding graph operation from the context of sorting by reversals [31] correspond to circle graphs and the (looped) local complementation operation, respectively, of the theory of circuit partitions of 4-regular multigraphs. This leads to a reformulation of the Hannenhalli-Pevzner theorem [27] to delta-matroids. We remark however that the theory of circuit partitions of 4-regular multigraphs is too broad to fully cover here, and so various references are provided in the paper for more information.

It has been shown that various other research topics can also be fit into the theory of circuit partitions of 4-regular multigraphs: examples include the theory of ribbon graphs (or embedded graphs) [8, 20] and the theory of gene assembly in ciliates [17]. As such all these research topics arising from different contexts turn out to be strongly linked, and results from one research topic can often be carried over to another. Indeed, it is not surprising that the Hannenhalli-Pevzner theorem has been independently discovered in the context of gene assembly in ciliates [23, 16].

This paper is organized as follows. In Section 2 we recall sorting by reversals and in Section 3 we associate a 4-regular multigraph and a pair of circuit partitions to a pair of chromosomes (where one is obtainable from the other by reversals). This leads to a reformulation of a known inequality of the reversal distance in terms of 4-regular multigraphs. In Section 4 we associate a circle graph to the two circuit partitions and we show that the adjacency matrix representation of this circle graph reveals essential information regarding the reversal distance. We recall local complementation in Section 5 and the Hannenhalli-Pevzner theorem in Section 6, and reformulate the Hannenhalli-Pevzner theorem in terms of delta-matroids in Section 8. Before discussing delta-matroids, we also recall sorting by DCJ operations in Section 7. We discuss the close connection of sorting by reversals and gene assembly in ciliates in Section 9. Finally, a discussion is given in Section 10.

## 2. Sorting by reversals

In this section we briefly and informally recall notions concerning sorting by reversals. See, e.g., the text books [34, 25] for a more formal and extensive treatment.

During the evolution of species, various types of modifications of the genome may occur. One such modification is the inversion (i.e., rotation by 180 degrees)

1	-6	7	4	-2	-5	3
---	----	---	---	----	----	---

chromosome of some species *A*

1	2	3	4	5	6	7
---	---	---	---	---	---	---

chromosome of some species *B*

Figure 1: Difference of two chromosomes of species *A* and *B*. The difference is described with respect to the ordering of the chromosome segments of *B*.

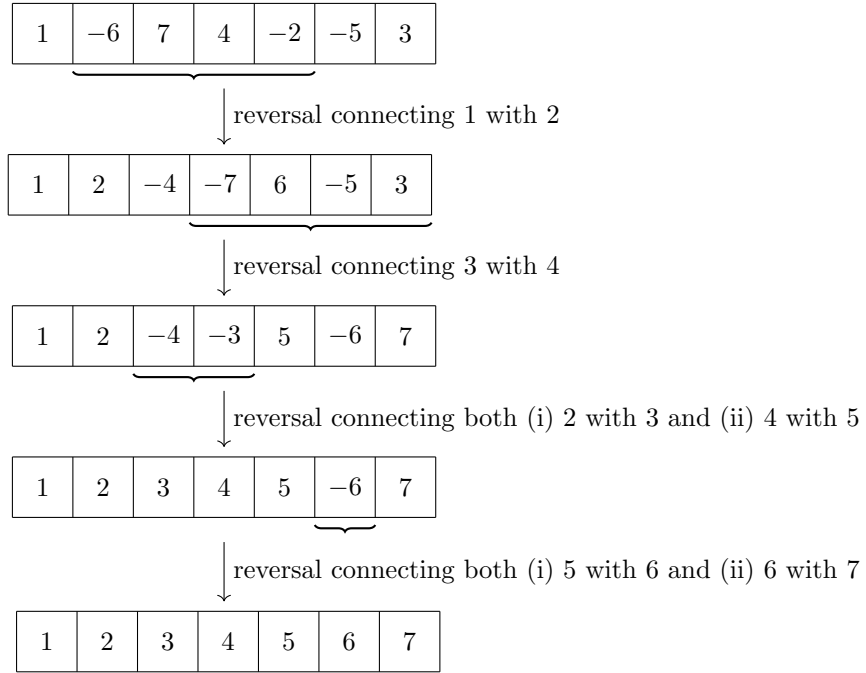


Figure 2: Optimal sorting of the signed permutation of Figure 1 by reversals.

of part of a chromosome, called a *reversal*. In Section 7, we recall that this inversion is the result of a so-called double-cut-and-join operation. The *reversal distance* between two given chromosomes is the minimal number of reversals needed to transform one into the other, and it is a measure of the evolutionary distance between the two species. Figure 1 shows two toy chromosomes which have seven segments in common, but their relative positions and orientations differ (by, e.g.,  $-5$  we mean segment 5 in inverted orientation, i.e., rotated by 180 degrees). Figure 2 shows that the reversal distance between the chromosomes of Figure 1 is at most four.

We can concisely describe the chromosomes of Figure 1 by the sequences  $(1, -6, 7, 4, -2, -5, 3)$  and  $(1, 2, 3, 4, 5, 6, 7)$ . These sequences are called *signed permutations* in the literature, and we will adopt this convention here (although we won't treat them as permutations in this paper). A signed permutation of the form  $(1, 2, \dots, n)$  is called the *identity permutation*. The *reversal distance* of a single signed permutation  $\pi$ , denoted by  $d_r(\pi)$ , is the minimal number of reversals needed to transform  $\pi$  into the identity permutation. Thus, for  $\pi = (1, -6, 7, 4, -2, -5, 3)$  of Figure 1, we have  $d_r(\pi) \leq 4$  by Figure 2.

Viewing the chromosome of species  $B$  of Figure 1 as the “sorted” chromosome, the transformation using reversals of the chromosome of species  $A$  to the chromosome of species  $B$  is called *sorting by reversals*.

**Remark 1.** For any signed permutation  $\pi = (\pi_1, \pi_2, \dots, \pi_n)$ , the signed permutation  $\bar{\pi} = (-\pi_n, -\pi_{n-1}, \dots, -\pi_1)$  represents the same chromosome (just considered 180 degrees rotated). The reversal distance of  $\pi$  and  $\bar{\pi}$  may however differ (this difference is, of course, at most one because  $\pi$  can be turned into  $\bar{\pi}$  using one “full” reversal). Therefore, one could argue that a better notion of the reversal distance of  $\pi$  would be the minimum value of the reversal distances of both  $\pi$  and  $\bar{\pi}$ . Equivalently, one could view both  $(1, 2, \dots, n)$  and  $(-n, -(n-1), \dots, -1)$  as identity permutations. We revisit this issue in Sections 3 and 7.

### 3. Four-regular multigraphs

In this paper, *graphs* are allowed to have loops but not multiple edges, and *multigraphs* are allowed to have both loops and multiple edges. We denote the sets of vertices and edges of a (multi)graph  $G$  by  $V(G)$  and  $E(G)$ , respectively. For graphs, each edge  $e \in E(G)$  is either of the form  $\{v\}$  (i.e.,  $v$  is a looped vertex) or of the form  $\{v_1, v_2\}$  (i.e., there is an edge between  $v_1$  and  $v_2$ ). A vertex  $v$  is said to be *isolated* if no edge is incident to it (in particular,  $v$  is not looped). A *4-regular multigraph* is a multigraph where each vertex has degree 4, a loop counting as two.

A standard tool for the calculation of the reversal distance is the so-called breakpoint graph of a signed permutation. In this section we instead assign a 4-regular multigraph and two of its circuit partitions to a signed permutation. The main reason for considering this graph instead of the breakpoint graph is that in this way we can use the vast amount of literature concerning the theory of

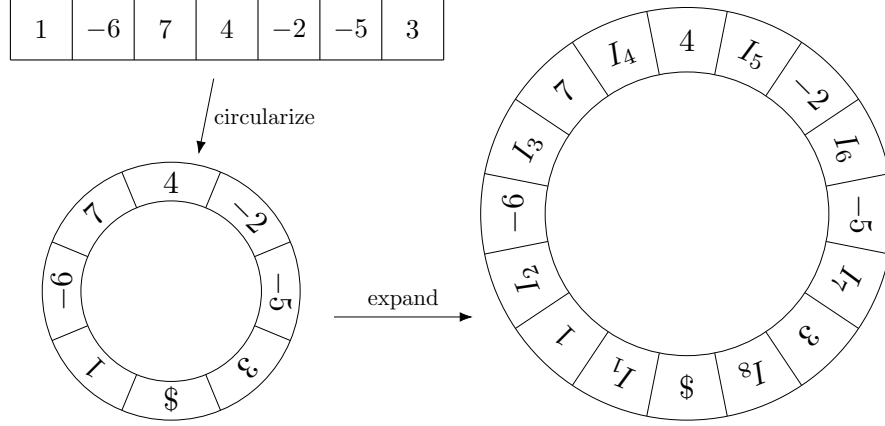


Figure 3: Circularizing and expanding the signed permutation of Figure 1.

circuit partitions of 4-regular multigraphs, which began with the seminal paper of Kotzig [33], and extended, e.g., in [1, 26]. We also note that a drawback of using the breakpoint graph is that the identity of the vertices is important, while the theory of circuit partitions of 4-regular multigraphs is independent of the identity of the vertices.

We now describe the construction of the 4-regular multigraph (along with two circuit partitions). As usual, the boundaries between adjacent segments of (the chromosomal depiction of) a signed permutation  $\pi$  are called *breakpoints*. Because reversals can also be applied on endpoints of a chromosome, we treat the endpoints of a signed permutation as breakpoints as well. We do this by circularizing the signed permutation, see Figure 3. Note, however, that the location of the endpoints is important. Indeed, e.g., the signed permutation  $(1, 2, \dots, n)$ , with  $n \geq 2$ , has a different reversal distance than any of its proper conjugations (i.e., the signed permutations  $(i, i+1, \dots, n, 1, 2, \dots, i-1)$  for  $i \in \{2, \dots, n\}$ ). Therefore, we have anchored the two endpoints to a new segment, which is denoted by  $\$$ . This corresponds to the usual procedure of *framing* the signed permutation in the theory of sorting by reversals, see, e.g., [3]. The next step, which we call here the *expand* step, is to insert an intermediate segment  $I_i$  between each two adjacent segments, see again Figure 3.

Next, we represent the circularized and expanded signed permutation by a digraph  $D_\pi$ . In  $D_\pi$ , each breakpoint is represented by a vertex and each segment is represented by an arrow. The arrow is labeled by the segment  $x$  it represents and goes from the left-hand breakpoint of  $x$  to the right-hand breakpoint of  $x$ , see Figure 4. Moreover, the boundaries/breakpoints of the original segments  $i$  and  $i+1$  that coincide after the sorting procedure are given a common vertex label  $v_i$ . For example, the right-hand side breakpoint of segment 3 and the left-hand side breakpoint of segment 4 are given a common vertex label  $v_3$ , see again Figure 4. Notice that the orientation is important here: the

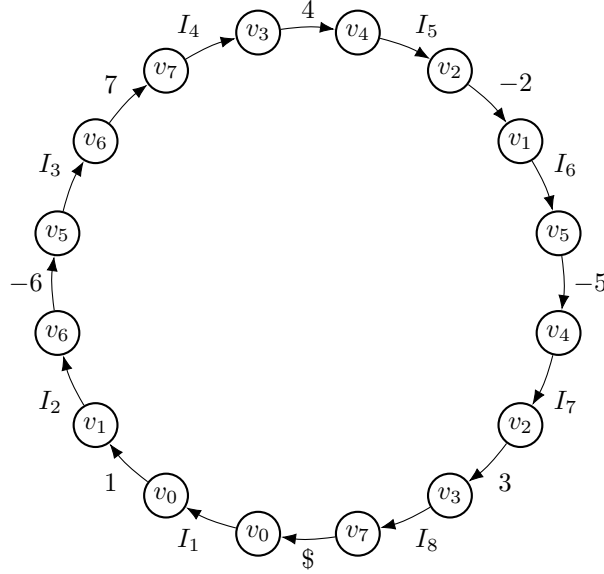


Figure 4: Digraph  $D_\pi$  representing the circularized and expanded signed permutation of Figure 3.

right-hand side breakpoints of segment 1 and the *right*-hand side breakpoint of segment  $-2$  are given a common vertex label because segment  $-2$  is segment 2 in inverted orientation (i.e., rotated by 180 degrees). Since segment  $\$$  represents the endpoints, for the arrow corresponding to  $\$$ , the head vertex is labeled by  $v_0$  and the tail vertex is labeled by  $v_n$ .

From the digraph  $D_\pi$  of Figure 4, we construct a 4-regular multigraph, denoted by  $G_\pi$ , by turning each arrow into an undirected edge, removing the signs from the edge labels, and finally merging each two vertices with the same label, see Figure 5.

Let  $G$  be a multigraph and let  $l$  be the number of connected components of  $G$ . A (unoriented) *circuit* of  $G$  is a closed walk, without distinguished orientation or starting vertex, allowing repetitions of vertices but not of edges. A *circuit partition* of  $G$  is a set  $P$  of circuits of  $G$  such that each edge of  $G$  is in exactly one circuit of  $P$ . Note that  $|P| \geq l$ . If  $|P| = l$ , then we say that  $P$  is an *Euler system* of  $G$ . Note that an Euler system contains an Eulerian circuit (i.e., a circuit visiting each edge exactly once) for each connected component of  $G$ . In particular, if  $G$  is connected, then an Euler system is a singleton containing an Eulerian circuit of  $G$ .

As illustrated in Figure 6, the circularized and expanded signed permutation of Figure 3 belongs to a particular circuit partition  $P_A$  of the 4-regular multigraph  $G_\pi$  of Figure 5 (this can also be verified by comparing Figure 6 with Figure 4). In this way,  $P_A$  is the circuit partition belonging to the chromosome

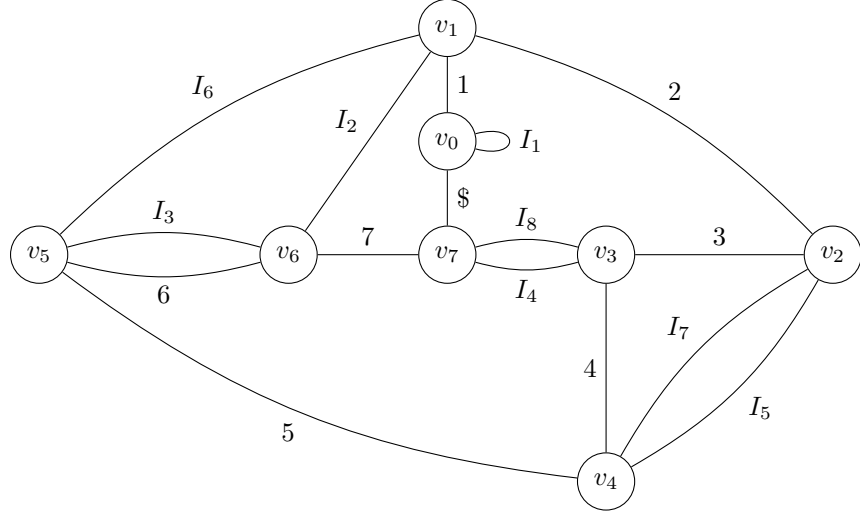


Figure 5: The 4-regular multigraph  $G_\pi$  of the running example.

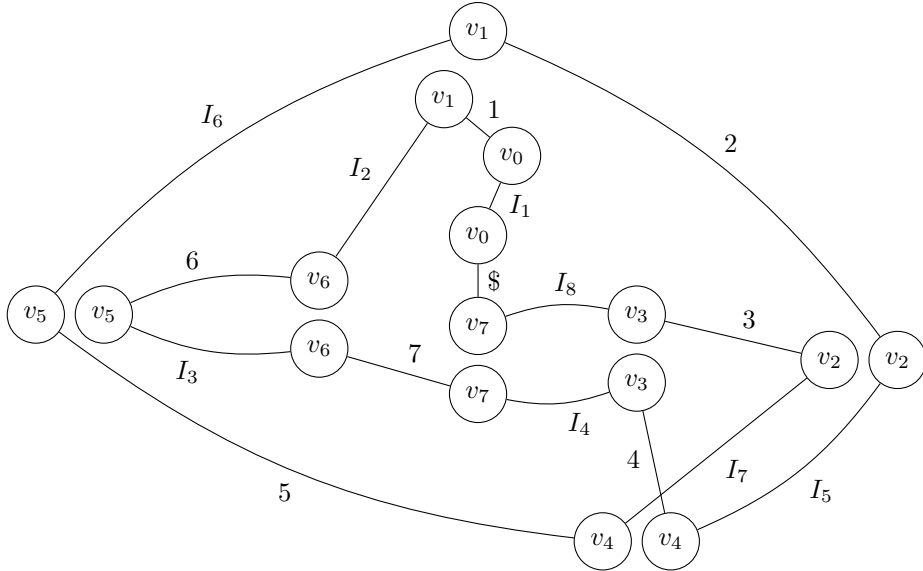


Figure 6: Circuit partition  $P_A$  of  $G_\pi$  representing the chromosome of species  $A$  of Figure 1.



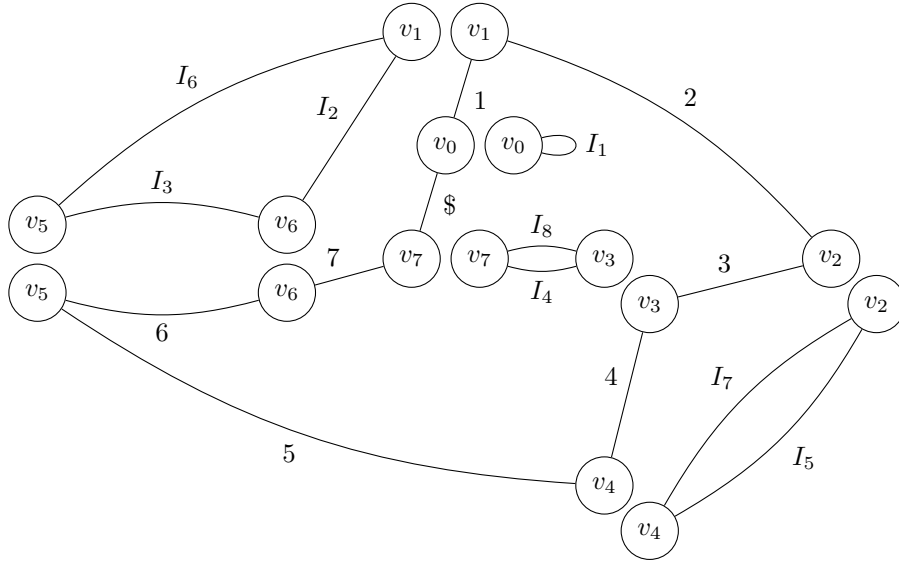


Figure 7: Circuit partition  $P_B$  of  $G_\pi$  containing a circuit  $C = (1, 2, \dots, 7, \$)$  that represents the chromosome of species  $B$  of Figure 1.

of species  $A$ . Notice that  $P_A$  is an Euler system, in fact, since  $G_\pi$  is connected for every signed permutation  $\pi$ ,  $P_A$  contains an Eulerian circuit of  $G_\pi$ . Another circuit partition  $P_B$  is illustrated in Figure 7. It is the unique circuit partition that includes the circuit  $C = (1, 2, \dots, 7, \$)$ . As such,  $P_B$  is the circuit partition belonging to the chromosome of species  $B$ . Notice that, besides  $C$ ,  $P_B$  contains four other circuits in this example. Each of these four circuits consists of intermediate segments (recall that these are the segments of the form  $I_i$  for some  $i$ ).

We remark that Figures 6 and 7 (corresponding to  $P_A$  and  $P_B$ , respectively) can be obtained from  $G_\pi$  by “splitting” each vertex in an appropriate way. This splitting is considered in [24] in the context of gene assembly in ciliates (we recall gene assembly in ciliates in Section 9).

While we do not recall the notion of a cycles of a signed permutation (see, e.g., [27, 3]), we mention that it is easy to verify that these cycles correspond one-to-one to circuits of intermediate segments of the circuit partition  $P_B$ . By using 4-regular multigraphs, we have given these cycles a more “physical” interpretation, cf. Figure 7.

Let  $c(\pi)$  be the number of cycles of a signed permutation  $\pi$ . The following result is well known (in fact, this result has been extended into an equality in [27]).

**Theorem 2 ([2]).** *Let  $\pi$  be a signed permutation with  $n$  elements. Then  $d_r(\pi) \geq n + 1 - c(\pi)$ .*

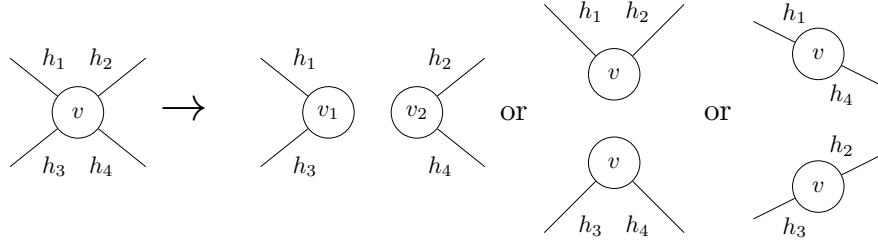


Figure 8: The three possible routes a circuit partition can take.

The proof idea of Theorem 2 is to show that (1) if  $\pi$  is the identity permutation, then  $c(\pi) = n + 1$  and (2) if  $\pi'$  is obtained from  $\pi$  by applying a single reversal, then  $c(\pi') - c(\pi) \in \{-1, 0, 1\}$ .

We remark that the inequality of Theorem 2 usually takes the form  $d_r(\pi) \geq n' - c(\pi)$ , where  $n'$  is the number of segments of the framed/anchored signed permutation and is natural when segment  $\$$  is instead denoted by  $n' = n + 1$ .

For our running example we see by Figure 7 that  $c(\pi) = 4$ . Thus  $d_r(\pi) \geq 7 + 1 - 4 = 4$ . We have seen in Section 2, and in particular Figure 2, that  $d_r(\pi) \leq 4$ . Consequently,  $d_r(\pi) = 4$ .

Notice that a circuit partition takes, for each vertex of the 4-regular multigraph, one of three possible routes, see Figure 8. Since care must be taken in the case of loops, the four  $h_i$ 's are not edges but actually “half-edges”, where two half-edges form an edge. If circuit partitions  $P_1$  and  $P_2$  take a different route at each vertex of  $G$ , then we say that  $P_1$  and  $P_2$  are *supplementary*. Notice that  $P_A$  and  $P_B$  are supplementary circuit partitions.

For a given circuit partition  $P$  and vertex  $v$  of  $G$ , let  $P'$  and  $P''$  be the circuit partitions obtained from  $P$  by changing the route of  $P$  at vertex  $v$ . It is well known that the cardinalities of two of  $\{P, P', P''\}$  are equal, to say  $k$ , and the third is of cardinality  $k + 1$ .

**Remark 3.** *In terms of 4-regular multigraphs, the issue discussed in Remark 1 translates to the question of whether the anchor should be  $\$$  or  $-\$$  (in other words,  $\$$  in inverted orientation). We assume the former, but the latter anchor is equally valid and may sometimes obtain a reversal distance that is one smaller. We revisit this issue in Section 7.*

#### 4. Circle graphs

To study the effect of sequences of reversals, we turn to circle graphs. Let us fix two supplementary circuit partitions  $P_1$  and  $P_2$  of a 4-regular multigraph  $G$ , where  $P_1$  is an Euler system.

A vertex  $v$  of  $G$  is called *oriented* for  $P_1$  with respect to  $P_2$  if the circuit partition  $P'$  obtained from  $P_1$  by changing the route of  $P_1$  at  $v$  to coincide with the route of  $P_2$  at  $v$ , is an Euler system. We say that vertex  $v$  of  $G_\pi$  is *oriented*

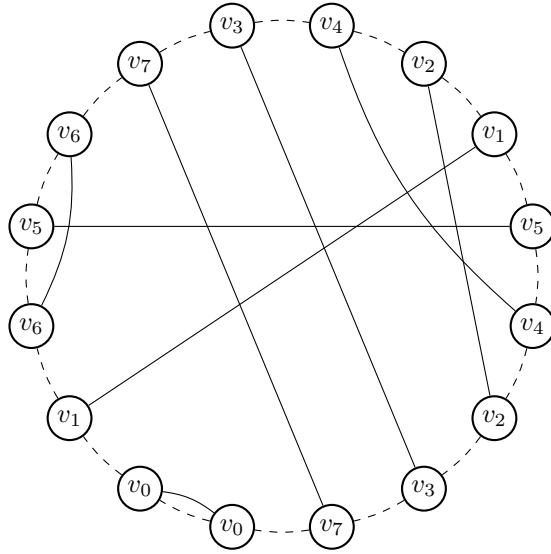


Figure 9: Chord diagram corresponding to the Eulerian circuit depicted in Figure 6.

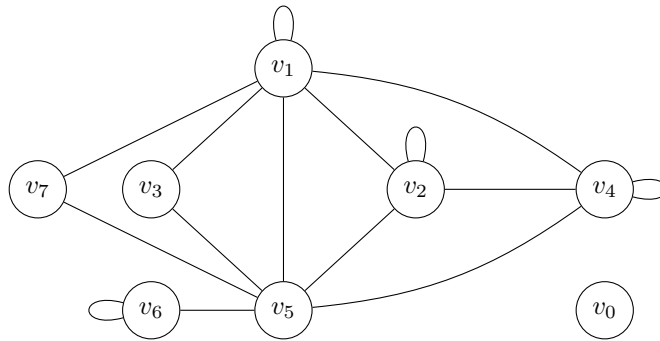


Figure 10: Circle graph  $H_\pi$  for the signed permutation  $\pi$  of the running example.

for a signed permutation  $\pi$  if  $v$  is oriented for  $P_A$  with respect to  $P_B$ . Thus  $\{v_1, v_2, v_4, v_6\}$  is the set of oriented vertices of our running example.

To construct the circle graph, we assume first, for convenience, that  $G$  is connected, i.e.,  $P_1$  contains a single Eulerian circuit  $C$ . We draw  $C$  as a circle and connect each two vertices with the same label by a chord to obtain a *chord diagram*. See Figure 9 for the chord diagram of the Eulerian circuit of Figure 6. We construct a (looped) *circle graph*  $H$  for  $G$  with respect to  $P_1$  and  $P_2$  as follows. The set  $V(H) = V(G)$  and two distinct vertices of  $H$  are adjacent when the corresponding two chords intersect in the chord diagram. Finally, a loop is added for each vertex that is oriented for  $P_1$  with respect to  $P_2$ . In the general case where  $G$  is not necessarily connected, the circle graph for  $G$  is the union of the circle graphs of each connected component of  $G$ . The *circle graph*  $H_\pi$  for signed permutation  $\pi$  is the circle graph of  $G_\pi$  with respect to  $P_A$  and  $P_B$ . The circle graph for our running example is depicted in Figure 10.

**Remark 4.** *We remark that a circle graph is called an “overlap graph” in the literature of sorting by reversals. However, we use here the term circle graph because circle graph is the usual name for this notion in mathematics. Also, the vertices of an overlap graph in the literature on sorting by reversals are often decorated by white or black labels instead of loops (black labels correspond to loops). The rest of this section shows why using loops instead of vertex colors is very useful when going to other combinatorial structures and matrices.*

We now recall the well-known notion of an adjacency matrix of a graph. First, the rows and columns of the matrices we consider in this paper are not ordered, but are instead indexed by finite sets  $X$  and  $Y$ , respectively. We call such matrices  $X \times Y$ -matrices. Note that the usual notions of rank and nullity of such a matrix  $A$  are defined — they are denoted by  $r(A)$  and  $n(A)$ , respectively. The *adjacency matrix* of a graph  $G$ , denoted by  $A(G)$ , is the  $V(G) \times V(G)$ -matrix over the binary field  $GF(2)$  where for  $v, v' \in V(G)$ , the entry indexed by  $(v, v')$  is 1 if and only if  $v$  and  $v'$  are adjacent (a vertex  $v$  is considered adjacent to itself precisely when  $v$  has a loop).

The adjacency matrix  $A(H_\pi)$  of the circle graph  $H_\pi$  of Figure 10 is as follows:

$$A(H_\pi) = \begin{matrix} & v_0 & v_1 & v_2 & v_3 & v_4 & v_5 & v_6 & v_7 \\ \begin{matrix} v_0 \\ v_1 \\ v_2 \\ v_3 \\ v_4 \\ v_5 \\ v_6 \\ v_7 \end{matrix} & \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 & 1 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 & 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 & 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 1 & 1 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 \end{pmatrix} \end{matrix}.$$

We now recall the following result from [36, Theorem 4].

**Theorem 5 ([36]).** *Let  $G$  be a 4-regular multigraph with  $l$  connected components and let  $P_1$  and  $P_2$  be supplementary circuit partitions of  $G$  with  $P_1$  an Euler system. Let  $H$  be the circle graph for  $G$  with respect to  $P_1$  and  $P_2$ . Then  $n(A(H)) = |P_2| - c$ .*

We have the following corollary to Theorem 5, which considers the case where  $H$  is of the form  $H_\pi$ .

**Corollary 6.** *Let  $\pi$  be a signed permutation. Then  $n(A(H_\pi)) = c(\pi)$ .*

PROOF. By Theorem 5,  $n(A(H_\pi)) = |P_B| - |P_A| = |P_B| - 1$  since  $P_A$  contains an Eulerian circuit of  $G_\pi$ . Recall from Section 3 that  $c(\pi)$  is the number of circuits of  $P_B$  excluding the circuit  $(1, \dots, n, \$)$ . Thus  $c(\pi) = |P_B| - 1$ .  $\square$

Corollary 6 illustrates the usefulness of using loops instead of vertex colors for circle graphs (cf. Remark 4).

For our running example, we have that  $c(\pi) = |P_B| - 1 = 5 - 1 = 4$ , so  $n(A(H_\pi)) = c(\pi) = 4$ .

Using Corollary 6, we can translate the inequality of Theorem 2 as follows.

**Theorem 7.** *Let  $\pi$  be a signed permutation. Then  $d_r(\pi) \geq r(A(H_\pi))$ .*

PROOF. Let  $\pi$  have  $n$  elements. By Theorem 2 and Corollary 6,  $d_r(\pi) \geq n + 1 - c(\pi) = n + 1 - n(A(H_\pi))$ . The result follows by observing that  $|V(H_\pi)| = n + 1$ .  $\square$

**Lemma 8.** *Let  $\pi$  be a signed permutation. Then  $\pi$  is the identity permutation if and only if  $r(A(H_\pi)) = 0$ .*

PROOF. Let  $\pi$  have  $n$  elements. Note that  $\pi$  is the identity permutation if and only if each intermediate segment forms a circuit of length 1 if and only if  $c(\pi) = n + 1$ . By Corollary 6, this is equivalent to  $n(A(H_\pi)) = n + 1$  and therefore equivalent to  $r(A(H_\pi)) = 0$ .  $\square$

Note that  $r(A(H)) = 0$  simply means that  $H$  contains no edges (i.e., consists of only isolated vertices).

## 5. Local complementation

In order to study the effect of reversals on circle graphs, we recall the following graph notions. For a graph  $H$  and vertex  $v$ , the *neighborhood* of  $v$  in  $H$ , denoted by  $N_H(v)$ , is  $\{v' \in V(G) \mid \{v, v'\} \in E(H), v' \neq v\}$ .

**Definition 9.** Let  $H$  be a graph and  $v$  a looped vertex of  $H$ . The *local complement* of  $H$  at  $v$ , denoted by  $H * v$ , is the graph obtained from  $H$  by complementing the subgraph induced by  $N_H(v)$ .

In other words, for all  $p \subseteq V(H)$  with  $|p| \in \{1, 2\}$ , we have  $p \in E(H * v)$  if and only if either (1)  $p \notin E(H)$  and  $p \subseteq N_H(v)$  or (2)  $p \in E(H)$  and  $p \not\subseteq N_H(v)$ .

Moreover, we denote by  $H *_c v$  the graph obtained from  $H * v$  by removing all edges incident to  $v$  (including the loop on  $v$ ). Thus  $v$  is an isolated vertex of  $H *_c v$ . Equivalently,  $H *_c v$  complements the subgraph induced by the “closed neighborhood”  $\{v' \in V(H) \mid \{v, v'\} \in E(H)\}$ . Finally, we denote by  $H|v$  the graph obtained from  $H *_c v$  by removing the isolated vertex  $v$ . Note that each of these three operations ( $*v$ ,  $*_c v$  and  $|v$ ) are only allowed on looped vertices  $v$ ; we say that such an operation is *applicable* to  $H$  if  $v$  is a looped vertex of  $H$ .

The interest of local complement for the topic of sorting by reversals is that it corresponds to applying some particular type of reversal [31].

**Theorem 10.** *Let  $\pi$  be a signed permutation and let  $\pi'$  be the signed permutation obtained from  $\pi$  by applying a reversal on the breakpoints corresponding to an oriented vertex  $v$ . Then  $H_{\pi'} = H_{\pi} *_c v$ .*

It turns out that local complementation has interesting effects on the underlying adjacency matrix (see, e.g., [15]).

**Lemma 11.** *Let  $H$  be a graph and  $v$  a looped vertex of  $H$ . Then  $n(A(H|v)) = n(A(H))$ . In other words,  $r(A(H *_c v)) = r(A(H)) - 1$ .*

While Lemma 11 can be proved directly, another way is to (1) observe that  $A(H|v)$  is obtained from  $A(H)$  by applying the Schur complementation matrix operation [35] on the submatrix induced by  $\{v\}$  and (2) recall from, e.g., [39] that Schur complementation preserves nullity. We remark that while  $A(H|v)$  is obtained from  $A(H)$  by applying Schur complementation,  $A(H * v)$  is obtained from  $A(H)$  by applying the principal pivot transform matrix operation, which is a partial matrix inversion operation, see, e.g., [38, 14].

Let  $H$  be a graph and let  $\sigma = (v_1, \dots, v_k)$  be a sequence of mutually distinct vertices of  $H$ . A sequence  $\varphi = *_c v_1 *_c v_2 \dots *_c v_k$  of  $*_c$  operations that is applicable to  $H$  (associativity of  $*_c$  is from left to right) is called an *lc-sequence* for  $H$ . We say that an lc-sequence is *full* if  $H\varphi$  contains only isolated vertices.

Since (1) by Lemma 11,  $*_c$  decreases rank by one and (2) a graph  $H$  contains only isolated vertices if and only if  $r(A(H)) = 0$ , we directly recover the following property observed in [21, Corollary 4] (see also [13, Section 6]).

**Corollary 12.** *Let  $H$  be a graph and let  $\varphi$  be an lc-sequence of length  $k$  for  $H$ . Then  $\varphi$  is full if and only if  $k = r(A(H))$ .*

Note that if a full lc-sequence  $\varphi$  exists for  $H_{\pi}$  with  $\pi$  a signed permutation, then by Theorem 7 and Lemma 8 we have that  $d_r(\pi) = r(A(H_{\pi}))$ . So, each full lc-sequence corresponds to an optimal sorting of  $\pi$ .

## 6. Hannenhalli-Pevzner theorem

We now recall the so-called Hannenhalli-Pevzner theorem which gives a precise criterion on arbitrary graphs  $H$  for the existence of a full lc-sequence. This theorem has been shown in [27] in terms of signed permutations  $\pi$ , but has later been extended to arbitrary graphs  $H$  (instead of essentially restricting to

circle graphs  $H_\pi$ ). Also, this result was shown independently in [23, 16] in the context of gene assembly in ciliates (we recall this topic in Section 9). We give here another proof of this result, closely following the reasoning of [3].

Let  $L$  be the set of looped vertices of a graph  $H$ . For all  $v \in V(H)$ , we denote  $N_H^l(v) = N_H(v) \cap L$  and  $N_H^{ul}(v) = N_H(v) \setminus L$ . Also,  $H$  is said to be *loopless* if  $L = \emptyset$ .

**Lemma 13.** *Let  $H$  be a connected graph and let  $v \in V(H)$  be looped. If  $H'$  is a loopless connected component of  $H|v$ , then both (1)  $V(H') \cap N_H^l(v) \neq \emptyset$  and (2)  $N_H^{ul}(v) \subseteq N_H^{ul}(w)$  and  $N_H^l(w) \subseteq N_H^l(v)$  for all  $w \in V(H') \cap N_H^l(v)$ .*

PROOF. Let  $H'$  be a loopless connected component of  $H|v$ . Since local complementation changes only edges between vertices of  $N_H(v)$ ,  $V(H') \cap N_H(v) \neq \emptyset$ . Because local complementation complements the loop status of each vertex of  $N_H(v)$  and  $H'$  is loopless,  $V(H') \cap N_H^l(v) = V(H') \cap N_H(v) \neq \emptyset$ .

Let  $w \in V(H') \cap N_H^l(v)$ .

Firstly, let  $x \in N_H^{ul}(v)$ . Then  $x$  is looped in  $H|v$ . Since  $H'$  is loopless,  $\{x, w\} \notin E(H|v)$ . Thus  $\{x, w\} \in E(H)$  (because  $x, w \in N_H(v)$ ). Consequently,  $x \in N_H^{ul}(w)$ . Thus  $N_H^{ul}(v) \subseteq N_H^{ul}(w)$ .

Secondly, let  $x \in N_H^l(w)$ . If  $x \notin N_H^l(v)$ , then  $x \in N_{H|v}^l(w)$  which contradicts the fact that  $w$  belongs to a loopless connected component of  $H|v$ . Thus  $N_H^l(w) \subseteq N_H^l(v)$ .  $\square$

For a vertex  $v$  of  $H$ , define  $s(v) = |N_H^{ul}(v)| - |N_H^l(v)|$ . Let  $\text{MS}(H)$  be the set of looped vertices  $v$  of  $H$  such that  $s(w) \leq s(v)$  for all  $w \in N_H^l(v)$ . Note that for any graph  $H$  with looped vertices,  $\text{MS}(H)$  is nonempty since it contains all looped vertices that are (globally) maximal with respect to function  $s$ .

**Lemma 14.** *Let  $H$  be a connected graph and  $v \in \text{MS}(H)$ . Then each loopless connected component of  $H|v$  consists of only an isolated vertex.*

PROOF. Let  $H'$  be a loopless connected component of  $H|v$ . Since  $v \in \text{MS}(H)$ , we have by Lemma 13,  $N_H^{ul}(v) = N_H^{ul}(w)$  and  $N_H^l(w) = N_H^l(v)$  for some  $w \in V(H')$ . Since  $v$  and  $w$  are moreover looped and adjacent, we have that  $w$ , and therefore  $H'$ , is an isolated vertex of  $H|v$ .  $\square$

By iteration of Lemma 14, we obtain the following.

**Theorem 15.** *Let  $H$  be a graph. Then there is a full lc-sequence for  $H$  if and only if each loopless connected component of  $H$  consists of only an isolated vertex.*

**Corollary 16 ([27]).** *Let  $\pi$  be a signed permutation. If each connected component of  $H_\pi$  has at least one looped vertex or consists of only an isolated vertex, then  $d_r(\pi) = n + 1 - c(\pi)$ .*

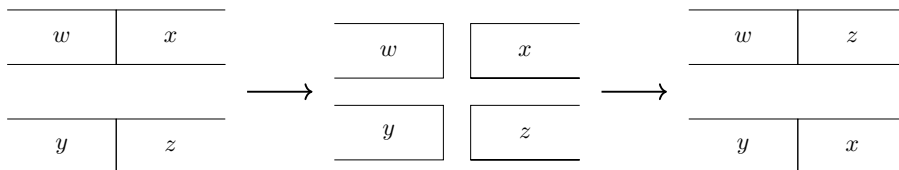


Figure 11: The DCJ operation.

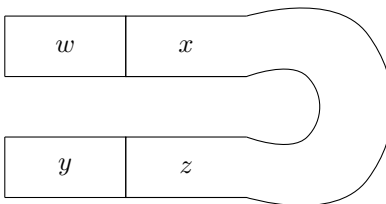


Figure 12: Alignment corresponding to a reversal.

## 7. DCJ operations and multiple chromosomes

A reversal is a special case of a *double-cut-and-join* (DCJ for short) operation, also called recombination in other contexts. A DCJ operation is depicted in Figure 11. A DCJ operation consists of three stages: first two distinct breakpoints align, then both breakpoints are cut, and finally the ends are glued back together as depicted in Figure 11. Since we consider endpoints to be breakpoints as well, any of  $w$ ,  $x$ ,  $y$  and  $z$  may be nonexistent. From the alignment of Figure 12 one observes that a reversal is a special case of a DCJ operation. DCJ operations are allowed to be intermolecular as well, and so sorting by DCJ operations may involve multiple chromosomes (for example a whole genome) and each chromosome may be linear or circular. The *DCJ distance* of two genomes  $g_A$  and  $g_B$ , denoted by  $d_{\text{DCJ}}(g_A, g_B)$ , is the minimal number of DCJ operations needed to transform one genome into the other. A toy example of two genomes  $g_A$  and  $g_B$  of species  $A$  and  $B$ , respectively, consisting of both linear and circular chromosomes is given in Figure 13.

If  $g_A$  and  $g_B$  contain only circular chromosomes, then the method of Section 3 to construct a 4-regular multigraph applies essentially unchanged — the only difference is that the circularization step is not done (and so no anchor  $\$$  is introduced) because the chromosomes are circular already. Thus, we directly apply the expand step to all chromosomes of  $g_A$  and then construct the 4-regular multigraph  $G$  as before. It is now a special case of [4, Theorem 1] that  $d_{\text{DCJ}}(g_A, g_B) = n - c$ , where  $n$  is the number of vertices of  $G$  (i.e.,  $n$  is the number of segments in  $g_A$  and  $g_B$ ) and  $c$  is the number of circuits containing intermediate segments of the circuit partition  $P_B$  belonging to  $g_B$ . The result



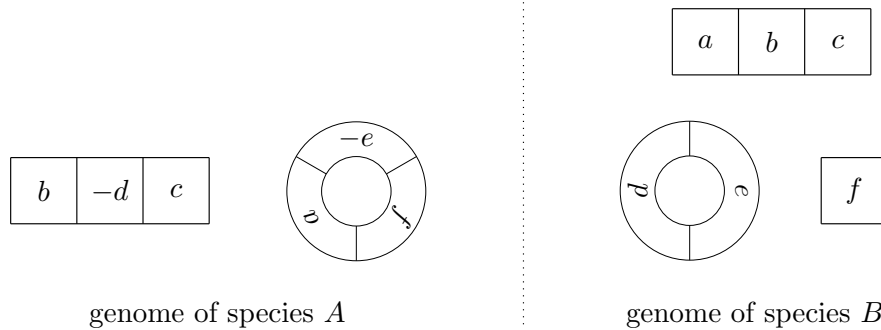


Figure 13: The genomes of species  $A$  and  $B$ .

in [4] is not stated in terminology of circuit partitions of a 4-regular multigraph, but instead in terms of a graph called an “adjacency graph”<sup>2</sup>. More precisely, in [4]  $c$  is defined as the number of cycles in the adjacency graph, and it can be readily verified that cycles in the adjacency graph correspond one-to-one to circuits of  $P_B$  containing intermediate segments.

While the construction of a 4-regular multigraph as outlined in Section 3 works when both genomes have only circular chromosomes, there are issues for genomes containing linear chromosomes such as in Figure 13. Indeed, the circularization of the linear chromosomes of  $g_B$  in Figure 13 leads to two anchors  $\$$ <sub>1</sub> and  $\$$ <sub>2</sub> which somehow need to be reconciled with the single anchor  $\$$  in the linear chromosome of  $g_A$ . Also, the issue discussed in Remarks 1 and 3 (which concerns the issue of whether to use  $\$$  or  $-\$$  as the anchor) is exasperated when there are several linear chromosomes. We leave it as an open problem to resolve this issue of the absence of a canonical 4-regular multigraph for two genomes. We mention that [4, Theorem 1] is able to calculate the DCJ distance in this general setting (i.e., with linear chromosomes). The formula takes the form  $d_{\text{DCJ}}(g_A, g_B) = n - (c + i/2)$ , where  $i$  is the number of connected components of the adjacency graph that are odd-length paths. Very roughly, one way to explain this formula in terms of 4-regular multigraphs is that each odd-length path corresponds to a side of a linear chromosome and an anchor (which increases  $n$  by one) can be introduced in such a way that both sides of a linear chromosome end up in different circuits (which increases  $c$  by two). So, the net effect of two sides of a linear chromosome is one, hence the contribution of  $i/2$ .

Recall that the construction of a circle graph from Section 4 requires two supplementary circuit partitions  $P_A$  and  $P_B$  of a 4-regular multigraph  $G$  where  $P_A$  is an Euler system. While in the theory of sorting by reversals  $P_A$  is always an Euler system (in fact,  $P_A$  contains an Eulerian circuit of  $G_\pi$ ), the circuit

<sup>2</sup>The notion of adjacency graph is not to be confused with the different notion of adjacency matrix as recalled earlier in this paper.

partition  $P_A$  belonging to a genome  $g_A$  is *not* necessarily a Euler system. In Section 8 we recall the notion of a delta-matroid which can function as a substitute for circle graphs in scenarios where circle graphs do not exist. In fact, as we will see, delta-matroids are useful even in scenarios where circle graphs do exist, such as in the theory of sorting by reversals.

## 8. Delta-matroids

It was casually remarked in the discussion of [28] that the Hannenhalli-Pevzner theorem might be generalizable to the more general setting of delta-matroids, which are combinatorial structures defined by Bouchet [5]. Indeed, we recall now how delta-matroids can be constructed from circuit partitions in 4-regular multigraphs and from circle graphs. In this way, various results concerning the theory of sorting by reversals can be viewed in this more general setting.

A *set system*  $D = (V, S)$  is an ordered pair, where  $V$  is a finite set, called the *ground set*, and  $S$  a set of subsets of  $V$ . For set systems  $D_1 = (V_1, S_1)$  and  $D_2 = (V_2, S_2)$  with disjoint ground sets, we define the *direct sum* of  $D_1$  and  $D_2$ , denoted by  $D_1 \oplus D_2$ , as the set system  $(V_1 \cup V_2, \{X_1 \cup X_2 \mid X_1 \in S_1, X_2 \in S_2\})$ . In this case we say that  $D_1$  and  $D_2$  are *summands* of  $D$ . A set system  $D$  is called *connected* if it is not the direct sum of two set systems with nonempty ground sets. Also,  $D = (V, S)$  is called *even* if the cardinalities of all sets in  $S$  are of equal parity. Let us denote symmetric difference by  $\Delta$ . A *delta-matroid*  $D = (V, S)$  is a set system where  $S$  is nonempty and, moreover, for all  $X, Y \in S$  and  $x \in X \Delta Y$ , there is an  $y \in X \Delta Y$  (possibly equal to  $x$ ) with  $X \Delta \{x, y\} \in S$  [5].

Let  $G$  be a 4-regular multigraph and let  $P_1$  and  $P_2$  be two supplementary circuit partitions of  $G$ . Denote by  $\mathcal{D}_G(P_1, P_2)$  the set system  $(V(G), S)$ , where for  $X \subseteq V(G)$  we have  $X \in S$  if and only if the circuit partition  $P$  obtained from  $P_1$  by changing, for each  $v \in X$ , the route of  $P_1$  at  $v$  to coincide with the route of  $P_2$  at  $v$ , is an Euler system.

The following result is stated in [5, Theorem 5.2] (see also [7, Theorem 5.2]).

**Theorem 17 ([5]).** *Let  $G$  be a 4-regular multigraph and let  $P_1$  and  $P_2$  be two supplementary circuit partitions of  $G$ . Then  $\mathcal{D}_G(P_1, P_2)$  is a delta-matroid.*

**Example 18.** Let  $G_\pi$  be from Figure 5,  $P_A = \{C\}$  be from Figure 6 and  $P_B$  be from Figure 7. Then  $\mathcal{D}_{G_\pi}(P_A, P_B) = (V(G), S)$ , where

$$\begin{aligned} S = & \{\emptyset, \{v_1\}, \{v_2\}, \{v_4\}, \{v_6\}, \{v_1, v_3\}, \{v_1, v_5\}, \{v_1, v_6\}, \{v_1, v_7\}, \{v_2, v_4\}, \{v_2, v_5\}, \\ & \{v_2, v_6\}, \{v_3, v_5\}, \{v_4, v_5\}, \{v_4, v_6\}, \{v_5, v_6\}, \{v_5, v_7\}, \{v_1, v_2, v_3\}, \dots\}. \end{aligned}$$

For a graph  $H$  and  $X \subseteq V(G)$ , we denote by  $H[X]$  the subgraph of  $H$  induced by  $X$  (i.e., all vertices outside  $X$  are removed including their incident edges). Also, denote by  $\mathcal{D}_H$  the set system  $(V(H), S)$ , where for  $X \subseteq V(H)$ ,  $X \in S$  if and only if the matrix  $A(H[X])$  is invertible. As usual, the empty matrix is invertible by convention.

**Theorem 19 (Theorem 4.1 of [7]).** *Let  $H$  be a graph. Then  $\mathcal{D}_H$  is a delta-matroid.*

The delta-matroids of the form  $\mathcal{D}_H$  for some graph  $H$  are called *binary normal delta-matroids* [7].

We now recall the close connection between delta-matroids of circle graphs and of circuit partitions of 4-regular multigraphs.

**Theorem 20 (Theorem 5.3 of [7]).** *Let  $H$  be the circle graph of a 4-regular multigraph  $G$  with respect to the supplementary circuit partitions  $P_1$  and  $P_2$  with  $P_1$  an Euler system. Then  $\mathcal{D}_H = \mathcal{D}_G(P_1, P_2)$ .*

Note that  $G_\pi$ ,  $P_A$  and  $P_B$  as used in Example 18 are also used to construct the circle graph  $H_\pi$  of Figure 10. So,  $\mathcal{D}_{G_\pi}(P_A, P_B)$  of Example 18 is equal to  $\mathcal{D}_{H_\pi}$ .

Notice, e.g., that the looped vertices of  $H$  precisely correspond to the singletons in  $\mathcal{D}_G(P_1, P_2)$ . Indeed, the adjacency matrix of a subgraph containing a single vertex  $v$  is invertible precisely when  $v$  is looped.

By Theorem 20,  $\mathcal{D}_G(P_1, P_2)$  corresponds to a circle graph if  $P_1$  an Euler system. Recall that in the case of the sorting of genomes by DCJ operations (see Section 7), the usual construction of a circle graph does not work since the input does not necessarily correspond to an Euler system. However, we can construct  $\mathcal{D}_G(P_1, P_2)$  and so delta-matroids allow one to work with combinatorial structures similar to circle graphs in cases (such as in Section 7) where a corresponding circle graph does not seem to exist. This is one important reason for considering delta-matroids. Another reason is that delta-matroids are often more easy to work with than circle graphs (even in the cases where circle graphs exist) since the operation of local complementation (on looped vertices) is much more simple in terms of delta-matroids. We will recall this now.

For a set system  $D = (V, S)$  and  $X \subseteq V$ , define the *twist* of  $D$  on  $X$ , denote by  $D * X$ , as the set system  $(V, S')$  with  $S' = \{X \Delta Y \mid Y \in S\}$ . Notice that  $(D * X) * X = D$ . Also, it is easy to verify that if  $D$  is a delta-matroid, then so is  $D * X$ .

**Theorem 21.** *Let  $H$  be a graph and  $v \in V(H)$  be looped. Then  $\mathcal{D}_{H*v} = \mathcal{D}_H * \{v\}$ .*

By Theorem 21,  $\mathcal{D}_{H*_c v}$  is obtained from  $\mathcal{D}_H * \{v\}$  by removing all sets containing  $v$ .

We can translate the notion of a full lc-sequence for graphs to the realm of delta-matroids. Let  $D = (V, S)$  be a set system. A sequence  $\sigma = (v_1, v_2, \dots, v_k)$  of mutually distinct elements of  $V$  is a *lc-sequence* for  $D$  if for all  $i \in \{0, \dots, k\}$ ,  $\{v_1, \dots, v_i\} \in S$ . Note that, in particular,  $\emptyset, \{v_1, \dots, v_k\} \in S$ . An lc-sequence  $\sigma = (v_1, v_2, \dots, v_k)$  for  $D$  is said to be *full* if  $\{v_1, \dots, v_k\} \in \max(S)$ , where  $\max(S)$  is the set of maximal elements of  $S$  with respect to inclusion.

**Lemma 22.** *Let  $H$  be a graph. Then  $\sigma = (v_1, v_2, \dots, v_k)$  is an lc-sequence for  $\mathcal{D}_H$  if and only if  $\varphi = *_c v_1 *_c v_2 \cdots *_c v_k$  is an lc-sequence for  $H$ . Moreover, in this case,  $\sigma$  is full if and only if  $\varphi$  is full.*

PROOF. We prove the first statement by induction on  $k$ . If  $k = 0$ , then  $\sigma$  and  $\varphi$  are empty sequences and so the result holds trivially. Assume that  $k > 0$  and that the result holds for  $k' = k - 1$ .

First, let  $\sigma$  be an lc-sequence for  $\mathcal{D}_H$ . In particular,  $\sigma' = (v_1, v_2, \dots, v_{k-1})$  is an lc-sequence for  $\mathcal{D}_H$ . By the induction hypothesis,  $\varphi' = *_c v_1 *_c v_2 \cdots *_c v_{k-1}$  is an lc-sequence for  $H$ . Assume to the contrary that  $\varphi$  is not an lc-sequence. Then  $v_k$  is not a looped vertex of  $H\varphi'$ . Thus  $\{v_k\}$  is not a set of  $\mathcal{D}_{H\varphi'}$ . By the sentence below Theorem 21,  $\{v_k\}$  is not a set of  $\mathcal{D}_H * \{v_1, v_2, \dots, v_{k-1}\}$ . In other words,  $\{v_1, v_2, \dots, v_k\}$  is not a set of  $\mathcal{D}_H$  — a contradiction.

Second, let  $\varphi$  be an lc-sequence for  $H$ . In particular,  $\varphi' = *_c v_1 *_c v_2 \cdots *_c v_{k-1}$  is an lc-sequence for  $H$ . By the induction hypothesis,  $\sigma' = (v_1, v_2, \dots, v_{k-1})$  is an lc-sequence for  $\mathcal{D}_H$ . It suffices now to show that  $\{v_1, v_2, \dots, v_k\}$  is a set of  $\mathcal{D}_H$ . By the sentence below Theorem 21,  $\{v_k\}$  is a set of  $\mathcal{D}_H * \{v_1, v_2, \dots, v_{k-1}\}$ . Thus  $\{v_1, v_2, \dots, v_k\}$  is a set of  $\mathcal{D}_H$ .

We now prove the second statement. By the strong principal minor theorem, see [32] and also [13, Lemma 12], each set in  $\max(S)$  is of cardinality  $r(A(H))$ . Thus lc-sequence  $\sigma$  is full if and only if  $k = r(A(H))$  if and only if  $\varphi$  is full (by Corollary 12).  $\square$

We are now ready to rephrase Theorem 15 in delta-matroid terminology as follows.

**Theorem 23.** *Let  $D$  be a binary normal delta-matroid. Then there is a full lc-sequence for  $D$  if and only if each even connected summand of  $D$  with nonempty ground set is of the form  $(\{v\}, \{\emptyset\})$  for some element  $v$  in the ground set of  $D$ .*

PROOF. Since  $D$  is a binary normal delta-matroid, we have that  $D = \mathcal{D}_H$  for some graph  $H$ . By Lemma 22, there is a full lc-sequence for  $D$  if and only if there is a full lc-sequence for  $H$ . Also observe that, if a graph  $H'$  consists of only an isolated vertex  $v$ , then  $\mathcal{D}_{H'} = (\{v\}, \{\emptyset\})$ . By Theorem 15, it suffices to recall that (1) a graph  $H'$  is loopless if and only if  $\mathcal{D}_{H'}$  is even (the if-direction is immediate and the only-if direction follows from the well-known fact that invertible zero-diagonal skew-symmetric matrices have even dimensions, see, e.g., [11]) and (2) for all  $X \subseteq V(G)$ , the subgraph of  $H$  induced by  $X$  is a (possibly empty) union of connected components of  $H$  if and only if there is a summand of  $\mathcal{D}_H$  with ground set  $X$  [9, Proposition 5].  $\square$

Theorem 23 does not hold for arbitrary delta-matroids. Indeed, the delta-matroid  $D = (V, S)$  with  $|V| = 3$  and  $S = \{X \subseteq V \mid |X| \neq 1\}$  is connected but not even and so the right-hand side of the equivalence of Theorem 23 trivially holds. However,  $D$  does not have a full lc-sequence since  $D$  does not contain any singletons. It would be interesting to see to which class of delta-matroids Theorem 23 can be generalized.

We briefly mention that delta-matroids have been generalized to multimatroids in [10]. In this general setting, delta-matroids translate to a class of multimatroids called 2-matroids. Therefore, this section could also have been phrased in the setting of 2-matroids. Another class of multimatroids, called

tight 3-matroids can also be associated to 4-regular multigraphs. See [12] for the case where the 4-regular multigraph is derived from the context of gene assembly in ciliates, which is a theory closely related to that of sorting by reversals, see Section 9. Although it is out of the scope of this paper, it would be interesting to study tight 3-matroids associated to 4-regular multigraphs from the context of sorting by reversals.

## 9. A closely related theory: gene assembly in ciliates

Gene assembly is a process taking place during sexual reproduction of unicellular organisms called ciliates. During this process, a nucleus, called the *micronucleus* (or MIC for short), is transformed into another, very different, nucleus, called *macronucleus* (or MAC for short). Each gene in the MAC is one block consisting of a number of consecutive<sup>3</sup> segments called MDSs. These MDSs also appear in the corresponding MIC gene, but they can appear there in (seemingly) arbitrary order and orientation with respect to the MAC gene and they are moreover separated by noncoding segments called IESs. A (toy) example of a gene consisting of seven MDSs  $M_1, \dots, M_7$  is given in Figures 14 and 15, where the IESs are denoted by  $I_1, \dots, I_8$  and if an MDS  $M_i$  in the MIC gene is in inverted orientation (i.e., rotated by 180 degrees) with respect to the MAC gene, then this is denoted by  $\overline{M_i}$  (this is the standard notation in this theory, and would of course be written with a minus sign, i.e.,  $-M_i$ , in the theory of sorting by reversals).

The postulated way in which a MIC gene is transformed into its MAC gene is through DCJ operations (called DNA recombination in this context), where MDSs that are not adjacent in the MIC gene but are adjacent in the MAC gene are aligned, cut and glued back such that the two MDSs become adjacent like they appear in the MAC. For example, segments  $w$  and  $z$  of Figure 11 may be MDSs  $M_i$  and  $M_{i+1}$ , respectively, and  $x$  and  $y$  IESs. In the *intramolecular* model, the application of DCJ operations is restricted in such a way that they cannot result in MDSs appearing in different molecules [22] (however, it is allowed to apply two DCJ operations simultaneously, when each of them separately would result in a split of the molecule).

Starting from the MIC gene, Figure 14, we can construct a 4-regular multigraph in a similar way as for the theory of sorting by reversals. However, unlike before, the left-hand side of the first MDS  $M_1$  and the right-hand side of the last MDS  $M_7$  are not considered breakpoints. Hence, in Figure 14, we merge adjacent segments  $I_1$  and  $M_1$  ( $M_7$  and  $I_4$ , respectively) in the MIC gene into a single segment called  $I_1M_1$  ( $M_7I_4$ , respectively). Moreover, recall that before we introduced an anchor segment  $\$$  during the circularization, because the endpoints of a chromosome are breakpoints too. In the case of MIC genes, there are no endpoints that are breakpoints and so we omit the anchor segment  $\$$ .

---

<sup>3</sup>Actually, the MDSs overlap slightly in the MAC (the overlapping regions are called point-ers), but this is not relevant for this paper.

$I_1$	$M_1$	$I_2$	$\overline{M_6}$	$I_3$	$M_7$	$I_4$	$M_4$	$I_5$	$\overline{M_2}$	$I_6$	$\overline{M_5}$	$I_7$	$M_3$	$I_8$
-------	-------	-------	------------------	-------	-------	-------	-------	-------	------------------	-------	------------------	-------	-------	-------

Figure 14: MIC form of a gene

	$M_1$	$M_2$	$M_3$	$M_4$	$M_5$	$M_6$	$M_7$	
--	-------	-------	-------	-------	-------	-------	-------	--

Figure 15: MAC form of a gene

Instead, during the circularization the rightmost segment  $I_8$  is merged with the leftmost segment  $I_1M_1$  in the MIC gene into a single segment called  $I_8; I_1M_1$ , see Figure 16. Finally, because the “signed permutation” of Figure 14 already contains intermediate segments (the IESs), we also do not have an expand step like in Figure 3. In this way, the intermediate segments are physical, while they are “virtual” in the theory of sorting by reversals.

The construction of a 4-regular multigraph from Figure 16 is now identical as in the theory of sorting by reversals. Due the similarity of the “signed permutations” of Figures 1 and 14, we see that the circle graph corresponding to Figure 16 is obtained from the circle graph of Figure 10 by removing vertices  $v_0$  and  $v_7$ . See [17] for more details on how the 4-regular multigraphs, circle graphs and delta-matroids can be used to study gene assembly in ciliates.

From the discussion above it is not surprising that the theory of sorting by reversals and the theory of gene assembly in ciliates partly overlap. For example, in both theories local complementation plays an important role — indeed, as we mentioned in Section 6 the Hannenhalli-Pevzner theorem was discovered independently in both theories. We remark that links between gene assembly in ciliates and sorting by reversals have also been appreciated in [18, 29, 30].

## 10. Discussion

We have shown that the theory of 4-regular multigraphs, including their accompanying combinatorial structures such as circle graphs and delta-matroids, can be applied to the topic of sorting by reversals. The fact that a notion such as local complementation has been rediscovered in the context of sorting by reversals signifies the importance of the theory of 4-regular multigraphs.

This paper may serve as an introduction to the theory of 4-regular multigraphs for the audience familiar with sorting by reversals. This paper has covered only very little of the extensive body of knowledge that the theory of 4-regular multigraphs provides and that can be applied to the topic of sorting by reversals. Indeed, we also mention for example the notion of *touch graph* [6, Section 6] (see also, e.g., [37]) of a circuit partition  $P$  that is likely to be useful for sorting by reversals. Indeed, while omitting details, the touch graph of a circuit partition  $P_B$  belonging to species  $B$  in the context of sorting by reversals

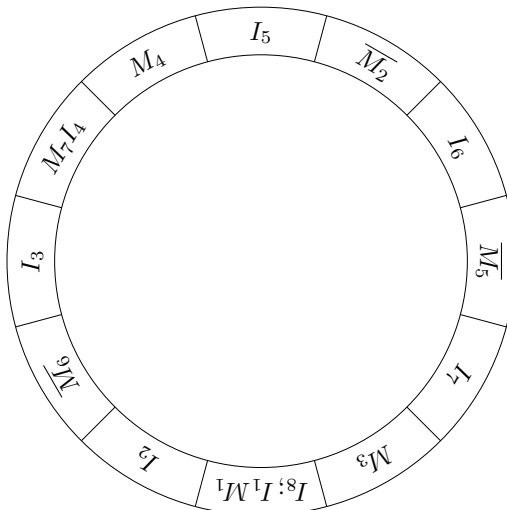


Figure 16: Circularization of the MIC gene of Figure 14.

(and also in the context of gene assembly in ciliates) is of a very special form, that of a star graph, which likely have interesting consequences. Indeed, while only implicitly stated in [19], this star graph property is key in the main result of [19] from the context of gene assembly in ciliates.

Since the theory of gene assembly in ciliates can also be fit into the theory of 4-regular multigraphs, we have also extended the links between the theories of gene assembly in ciliates and sorting by reversals as observed in [18, 29, 30].

Finally, we have formulated the open problem of using the theory of 4-regular multigraphs in the case of DCJ operations in the presence of linear chromosomes (see Section 7) and also the open problem of generalizing the Hannenhalli-Pevzner theorem from graphs to a suitable subclass of delta-matroids (see Section 8).

#### Acknowledgements

We thank anonymous referees for their helpful comments on an earlier version of this paper.

#### References

- [1] J. Abrham and A. Kotzig. Transformations of Euler tours. In M. Deza and I. G. Rosenberg, editors, *Combinatorics 79 Part I*, volume 8 of *Annals of Discrete Mathematics*, pages 65–69. Elsevier, 1980.
- [2] V. Bafna and P. A. Pevzner. Genome rearrangements and sorting by reversals. *SIAM Journal on Computing*, 25(2):272–289, 1996.

- [3] A. Bergeron. A very elementary presentation of the Hannenhalli-Pevzner theory. *Discrete Applied Mathematics*, 146(2):134–145, 2005.
- [4] A. Bergeron, J. Mixtacki, and J. Stoye. A unifying view of genome rearrangements. In P. Bucher and B. M. E. Moret, editors, *Proceedings of the 6th International Workshop on Algorithms in Bioinformatics (WABI 2006)*, volume 4175 of *Lecture Notes in Computer Science*, pages 163–173. Springer, 2006.
- [5] A. Bouchet. Greedy algorithm and symmetric matroids. *Mathematical Programming*, 38(2):147–159, 1987.
- [6] A. Bouchet. Isotropic systems. *European Journal of Combinatorics*, 8(3):231–244, 1987.
- [7] A. Bouchet. Representability of  $\Delta$ -matroids. In *Proceedings of the 6th Hungarian Colloquium of Combinatorics, Colloquia Mathematica Societatis János Bolyai*, volume 52, pages 167–182. North-Holland, 1987.
- [8] A. Bouchet. Maps and  $\Delta$ -matroids. *Discrete Mathematics*, 78(1-2):59–71, 1989.
- [9] A. Bouchet. Matroid connectivity and fundamental graphs. In *Proceedings of the 22nd Southeastern Conference on Combinatorics, Graph Theory, and Computing*, volume 85 of *Congressus Numerantium*, pages 81–88, 1991.
- [10] A. Bouchet. Multimatroids I. Coverings by independent sets. *SIAM Journal on Discrete Mathematics*, 10(4):626–646, 1997.
- [11] A. Bouchet and A. Duchamp. Representability of  $\Delta$ -matroids over  $GF(2)$ . *Linear Algebra and its Applications*, 146:67–78, 1991.
- [12] R. Brijder. Recombination faults in gene assembly in ciliates modeled using multimatroids. *Theoretical Computer Science*, 608:27–35, 2015.
- [13] R. Brijder and H. J. Hoogeboom. Maximal pivots on graphs with an application to gene assembly. *Discrete Applied Mathematics*, 158(18):1977–1985, 2010.
- [14] R. Brijder and H. J. Hoogeboom. The group structure of pivot and loop complementation on graphs and set systems. *European Journal of Combinatorics*, 32:1353–1367, 2011.
- [15] R. Brijder and H. J. Hoogeboom. Nullity invariance for pivot and the interlace polynomial. *Linear Algebra and its Applications*, 435:277–288, 2011.
- [16] R. Brijder and H. J. Hoogeboom. Binary symmetric matrix inversion through local complementation. *Fundamenta Informaticae*, 116(1-4):15–23, 2012.



- [17] R. Brijder and H. J. Hoogeboom. The algebra of gene assembly in ciliates. In N. Jonoska and M. Saito, editors, *Discrete and Topological Models in Molecular Biology*, Natural Computing Series, pages 289–307. Springer, 2014.
- [18] R. Brijder, H. J. Hoogeboom, and G. Rozenberg. Reducibility of gene patterns in ciliates using the breakpoint graph. *Theoretical Computer Science*, 356:26–45, 2006.
- [19] R. Brijder, H. J. Hoogeboom, and G. Rozenberg. Reduction graphs from overlap graphs for gene assembly in ciliates. *International Journal of Foundations of Computer Science*, 20:271–291, 2009.
- [20] C. Chun, I. Moffatt, S. Noble, and R. Rueckriemen. Matroids, delta-matroids and embedded graphs. [arXiv:1403.0920], 2014.
- [21] J. Cooper and J. Davis. Successful pressing sequences for a bicolored graph and binary matrices. *Linear Algebra and its Applications*, 490:162–173, 2016.
- [22] A. Ehrenfeucht, T. Harju, I. Petre, D. M. Prescott, and G. Rozenberg. *Computation in Living Cells – Gene Assembly in Ciliates*. Springer Verlag, 2004.
- [23] A. Ehrenfeucht, T. Harju, I. Petre, and G. Rozenberg. Characterizing the micronuclear gene patterns in ciliates. *Theory of Computing Systems*, 35:501–519, 2002.
- [24] A. Ehrenfeucht, T. Harju, and G. Rozenberg. Gene assembly through cyclic graph decomposition. *Theoretical Computer Science*, 281:325–349, 2002.
- [25] G. Fertin, A. Labarre, I. Rusu, E. Tannier, and S. Vialette. *Combinatorics of Genome Rearrangements*. MIT Press, 2009.
- [26] H. Fleischner, G. Sabidussi, and E. Wenger. Transforming eulerian trails. *Discrete Mathematics*, 109(1):103–116, 1992.
- [27] S. Hannenhalli and P. A. Pevzner. Transforming cabbage into turnip: Polynomial algorithm for sorting signed permutations by reversals. *Journal of the ACM*, 46(1):1–27, 1999.
- [28] T. Hartman and E. Verbin. Matrix tightness: A linear-algebraic framework for sorting by transpositions. In F. Crestani, P. Ferragina, and M. Sanderson, editors, *Proceedings of the 13th International Conference on String Processing and Information Retrieval (SPIRE 2006)*, volume 4209 of *Lecture Notes in Computer Science*, pages 279–290. Springer, 2006.
- [29] J. L. Herlin, A. Nelson, and M. Scheepers. Using ciliate operations to construct chromosome phylogenies. *Involve*, 9(1):1–26, 2016.

- [30] C. L. Jansen, M. Scheepers, S. L. Simon, and E. Tatum. Context directed reversals and the ciliate decryptome. [arXiv:1603.06149v3], 2016.
- [31] H. Kaplan, R. Shamir, and R. E. Tarjan. A faster and simpler algorithm for sorting signed permutations by reversals. *SIAM Journal on Computing*, 29(3):880–892, 1999.
- [32] V. Kodyalam, T. Lam, and R. Swan. Determinantal ideals, Pfaffian ideals, and the principal minor theorem. In *Noncommutative Rings, Group Rings, Diagram Algebras and Their Applications*, pages 35–60. American Mathematical Society, 2008.
- [33] A. Kotzig. Eulerian lines in finite 4-valent graphs and their transformations. In *Theory of graphs, Proceedings of the Colloquium, Tihany, Hungary, 1966*, pages 219–230. Academic Press, New York, 1968.
- [34] P. A. Pevzner. *Computational Molecular Biology: An Algorithmic Approach*. MIT Press, 2000.
- [35] J. Schur. Über Potenzreihen, die im Innern des Einheitskreises beschränkt sind. *Journal für die reine und angewandte Mathematik*, 147:205–232, 1917.
- [36] L. Traldi. Binary nullity, Euler circuits and interlace polynomials. *European Journal of Combinatorics*, 32(6):944–950, 2011.
- [37] L. Traldi. The transition matroid of a 4-regular graph: An introduction. *European Journal of Combinatorics*, 50:180 – 207, 2015.
- [38] M. Tsatsomeros. Principal pivot transforms: properties and applications. *Linear Algebra and its Applications*, 307(1-3):151–165, 2000.
- [39] F. Zhang, editor. *The Schur Complement and Its Applications*, volume 4 of *Numerical Methods and Algorithms*. Springer, 2005.