

# Genome Aliquoting Revisited

ROBERT WARREN<sup>1</sup> and DAVID SANKOFF<sup>2</sup>

## ABSTRACT

**We prove that the genome aliquoting problem, the problem of finding a recent polyploid ancestor of a genome, with breakpoint distance can be solved in polynomial time. We propose an aliquoting algorithm that is a 2-approximation for the genome aliquoting problem with double cut and join distance, improving upon the previous best solution to this problem, Feijão and Meidanis' 4-approximation algorithm.**

**Key words:** algorithms, evolution, genomic rearrangements.

## 1. INTRODUCTION

**C**OMPARING TWO GENOMES WITH DUPLICATED GENES IS DIFFICULT. None of the distances used to compare genomes today (e.g., breakpoint distance, reversal distance, double cut and join distance) handle duplicated genes. However, in the special case where all genes are duplicated the same number of times, there has been some success.

Informally, the *genome aliquoting problem* is the problem of finding a genome with one copy of every gene given a genome with exactly  $p$  copies of every gene such that the distance between the given and resulting genomes is minimized according to some distance metric. Thus, the genome aliquoting problem eliminates the duplicate genes allowing a genome to be compared with other genomes using an existing algorithm. Solving this problem will allow genomes which have undergone a recent polyploidization event, common in plants, to be compared.

There have been a number of solutions to the genome aliquoting problem where the genome has exactly two copies of every gene. This restricted version of the problem is called the *genome halving problem* and was first introduced in El-Mabrouk and Sankoff (2003). It was solved for reversal and translocation distance in El-Mabrouk and Sankoff (2003) and Alekseyev and Pevzner (2007), and for double cut and join distance in Warren and Sankoff (2009b) and Mixtacki (2008).

The *genome aliquoting problem* was introduced in Warren and Sankoff (2009a) along with a sketch of a heuristic algorithm for the problem under double cut and join distance. Feijão and Meidanis (2009) provided an exact solution under single cut or join distance which is also a 4-approximation algorithm under double cut and join distance. In this article, we provide an exact polynomial-time algorithm under breakpoint distance which is also a 2-approximation algorithm under double cut and join distance. Since our algorithm is similar to that presented in Warren and Sankoff (2009a), it also bounds that heuristic as a 2-approximation for double cut and join distance.

---

<sup>1</sup>School of Information Technology and Engineering and <sup>2</sup>Department of Mathematics and Statistics, University of Ottawa, Ottawa, Canada.

## 2. DUPLICATED GENOMES

The fundamental elements that we study are *genes*. We represent each gene as a pair of *extremities* such that a gene  $x$  is represented by its *head*  $\vec{x}$ , which corresponds to the 3' end of the gene, and its *tail*  $\check{x}$ , which corresponds to the 5' end of the gene. A *genome* is represented by a multiset of 2-multisets and 1-sets of extremities, called *adjacencies* and *telomeres* respectively, such that, for each gene, the genome contains exactly the same number of heads as tails. The functions  $\mathcal{G}(G)$  and  $\mathcal{E}(G)$  return the set of all genes or extremities respectively of a genome  $G$ . The collection of copies of the same gene are called a *gene family* with the *size* being the number of copies in the genome (i.e., the number of heads, or, alternatively, the number of tails, that appear in a genome).

Genomes with one or more gene families with size greater than one are challenging to manipulate. Not only must there be the same number of heads and tails, but to understand the layout of the genome, we must know for duplicated genes which head corresponds to which tail; otherwise, the layout of the genome is not unique. Similarly, to compute the distance of two genomes, we must know which copy in one genome matches which copy in another genome. Genomes where this information is known are called *ordered genomes*. In such genomes, we distinguish members of the same gene family by a subscript such that each head corresponds to the tail with the same subscript, e.g.,  $\vec{a}_1$  corresponds to  $\check{a}_1$ . Similarly, if we know the ordering between two genomes then each extremity in one genome must correspond to the extremity with the same subscript in the other genome, e.g.,  $\vec{a}_1$  in one genome must correspond to  $\vec{a}_1$  in the other genome. By default, we assume most genomes are ordered with themselves, but the challenge of comparing genomes with duplicated genes comes from finding an ordering between two genomes as different orderings change the distance. Thus, to compare genomes with duplicated genes, we must find an ordering between them that minimizes their distance.

For most situations, ordered genomes are all that is needed. However, ordered genomes can be difficult to use, e.g.,  $\{\vec{a}_1, \vec{b}_1\} \cap \{\vec{a}_2, \vec{b}_2\} = \emptyset$ , and yet, frequently, it is desirable to perceive them as equal. Since these problems tend to occur frequently in our work, we introduce the concept of an *unordered genome* where no effort is made to distinguish the genes. Let  $\tilde{G}$  be the unordered counterpart of an ordered genome  $G$ . Similarly, if  $\alpha$  is an element of  $G$ , then  $\tilde{\alpha}$  is an element of  $\tilde{G}$  and if  $S$  is a subset of  $G$  then  $\tilde{S}$  is a subset of  $\tilde{G}$ . Unordered extremities are distinguished from ordered extremities by the absence or presence respectively of a subscript. Transforming an ordered genome to an unordered genome is trivial but the reverse can be very difficult and is problem specific.

The key concept that helps solve the genome aliquoting problem is *perfection*:

**Definition 1.** We define perfection between two elements of a genome to be that the two elements are the same or are disjoint. For adjacencies we also have the special requirement that they are not of the form  $\{x, x\}$  for some extremity  $x$ . Since there are two types of elements, adjacencies and telomeres, two elements  $\alpha$  and  $\beta$  can be in one of the following three configurations:

- let  $\alpha$  and  $\beta$  be a pair of adjacencies. They are perfect if and only if they are disjoint or the same and neither are of the form  $\{x, x\}$  for some extremity  $x$ ;
- let  $\alpha$  and  $\beta$  be a pair of telomeres. Since all telomeres are either the same or disjoint all pairs of telomeres are perfect;
- let  $\alpha$  be an adjacency and  $\beta$  be a telomere, since they cannot be the same they must be disjoint in order to be perfect. Also,  $\alpha$  must not be of the form  $\{x, x\}$  for some extremity  $x$ ;

A genome, or any set of adjacencies and telomeres, is perfect if and only if all pairs of elements are perfect.

For example, given a genome  $G = \{\{\vec{a}_1\}, \{\vec{a}_1, \vec{b}_1\}, \{\vec{b}_1, \vec{d}_1\}, \{\vec{d}_1, \vec{e}_1\}, \{\vec{e}_1, \vec{b}_2\}, \{\vec{b}_2\}, \{\vec{a}_2\}, \{\vec{a}_2, \vec{c}_1\}, \{\vec{c}_1, \vec{c}_2\}, \{\vec{c}_2, \vec{d}_2\}, \{\vec{d}_2, \vec{e}_2\}, \{\vec{e}_2\}\}$  its corresponding unordered genome is  $\tilde{G} = \{\{\vec{a}\}, \{\vec{a}, \vec{b}\}, \{\vec{b}, \vec{d}\}, \{\vec{d}, \vec{e}\}, \{\vec{e}, \vec{b}\}, \{\vec{b}\}, \{\vec{a}\}, \{\vec{a}, \vec{c}\}, \{\vec{c}, \vec{c}\}, \{\vec{c}, \vec{d}\}, \{\vec{d}, \vec{e}\}, \{\vec{e}\}\}$ . In  $\tilde{G}$ , the adjacencies  $\{\vec{a}, \vec{b}\}$  and  $\{\vec{b}, \vec{d}\}$  are perfect because they are disjoint but  $\{\vec{a}, \vec{b}\}$  and  $\{\vec{a}, \vec{c}\}$  are not perfect because they are neither disjoint nor equal (they have  $\vec{a}$  in common but do not share the other extremity). The adjacency  $\{\vec{c}, \vec{c}\}$  can never be perfect because it is of the form  $\{x, x\}$  for some extremity  $x$ . The telomere  $\{\vec{a}\}$  is perfect when compared with both  $\{\vec{a}\}$  and  $\{\vec{e}\}$  since it is equal to the former and disjoint with the latter.  $\{\vec{e}\}$  is not perfect when compared with the adjacency  $\{\vec{e}, \vec{b}\}$  because they are not disjoint.

In biology, a *polyploid* is a genome with multiple copies of the same chromosome. A perfect unordered genome where all gene families are of the same size shares this property and, hence, is called a *polyploid*. Similarly, an ordered genome  $G$  whose corresponding unordered genome  $\tilde{G}$  is a polyploid is a *polyploid*.

For example, the ordered genome  $G = \{\{\vec{a}_1\}, \{\vec{a}_1, \vec{b}_1\}, \{\vec{b}_1\}, \{\vec{b}_2\}, \{\vec{a}_2, \vec{b}_2\}, \{\vec{a}_2\}\{\vec{e}_2\}, \{\vec{d}_2, \vec{e}_2\}, \{\vec{c}_2, \vec{d}_2\}, \{\vec{c}_2\}, \{\vec{e}_1\}, \{\vec{d}_1, \vec{e}_1\}, \{\vec{c}_1, \vec{d}_1\}, \{\vec{c}_1\}\}$  is a polyploid because its corresponding unordered genome  $\tilde{G} = \{\vec{a}\}, \{\vec{a}, \vec{b}\}, \{\vec{b}\}, \{\vec{b}\}, \{\vec{a}, \vec{b}\}, \{\vec{a}\}\{\vec{e}\}, \{\vec{d}, \vec{e}\}, \{\vec{c}, \vec{d}\}, \{\vec{c}\}, \{\vec{e}\}, \{\vec{d}, \vec{e}\}, \{\vec{c}, \vec{d}\}, \{\vec{c}\}$  is perfect and all gene families are of the same size (there are two genes in each family).

Our problem is as follows:

**Definition 2.** Given an ordered genome  $G$  with all gene families of size  $p$ , the genome aliquoting problem is to find a polyploid  $H$  with all gene families of size  $p$  ordered in relation to  $G$  such that the distance between  $G$  and  $H$  is minimal.

### 3. BREAKPOINT DISTANCE

A *breakpoint (BP)* is a difference in adjacencies or telomeres between two genomes, e.g., in  $G = \{\{\vec{a}_1\}, \{\vec{a}_1, \vec{b}_1\}, \{\vec{b}_1, \vec{d}_1\}, \{\vec{d}_1, \vec{e}_1\}, \{\vec{e}_1, \vec{b}_2\}, \{\vec{b}_2\}, \{\vec{a}_2\}, \{\vec{a}_2, \vec{c}_1\}, \{\vec{c}_1, \vec{c}_2\}, \{\vec{c}_2, \vec{d}_2\}, \{\vec{d}_2, \vec{e}_2\}, \{\vec{e}_2\}\}$  and  $H = \{\{\vec{a}_1\}, \{\vec{a}_1, \vec{b}_1\}, \{\vec{b}_1, \vec{c}_1\}, \{\vec{c}_1, \vec{d}_1\}, \{\vec{d}_1, \vec{e}_1\}, \{\vec{e}_1, \vec{a}_2\}, \{\vec{a}_2, \vec{b}_2\}, \{\vec{b}_2, \vec{c}_2\}, \{\vec{c}_2, \vec{d}_2\}, \{\vec{d}_2, \vec{e}_2\}, \{\vec{e}_2\}\}, \{\vec{b}_1, \vec{d}_1\}$  is a breakpoint in  $G$  since there is no equivalent element in  $H$  but  $\{\vec{a}_1\}$  is not a breakpoint because it is in both genomes. Extending this notion to entire genomes we arrive at the following definition:

**Definition 3.** The breakpoint distance, introduced in Sankoff and Blanchette (1997), between two genomes  $G$  and  $H$  is the number of breakpoints between  $G$  and  $H$ .

The breakpoint distance is easy to calculate. From Tannier et al. (2009), given two genomes  $G$  and  $H$ , let  $\mathcal{A}(G, H)$  be the set of adjacencies shared between  $G$  and  $H$  and let  $\mathcal{T}(G, H)$  be the set of telomeres then the *breakpoint distance* between  $G$  and  $H$  can be calculated using the following equation:

$$\mathcal{D}_{BP}(G, H) = |\mathcal{G}(G)| - |\mathcal{A}(G, H)| - \frac{|\mathcal{T}(G, H)|}{2} \quad (1)$$

For example, given two duplicated genomes  $G = \{\{\vec{a}_1\}, \{\vec{a}_1, \vec{b}_1\}, \{\vec{b}_1, \vec{d}_1\}, \{\vec{d}_1, \vec{e}_1\}, \{\vec{e}_1, \vec{b}_2\}, \{\vec{b}_2\}, \{\vec{a}_2\}\{\vec{a}_2, \vec{c}_1\}, \{\vec{c}_1, \vec{c}_2\}, \{\vec{c}_2, \vec{d}_2\}, \{\vec{d}_2, \vec{e}_2\}, \{\vec{e}_2\}\}$  and  $H = \{\{\vec{a}_1\}, \{\vec{a}_1, \vec{b}_1\}, \{\vec{b}_1, \vec{c}_1\}, \{\vec{c}_1, \vec{d}_1\}, \{\vec{d}_1, \vec{e}_1\}, \{\vec{e}_1, \vec{a}_2\}, \{\vec{a}_2, \vec{b}_2\}, \{\vec{b}_2, \vec{c}_2\}, \{\vec{c}_2, \vec{d}_2\}, \{\vec{d}_2, \vec{e}_2\}, \{\vec{e}_2\}\}$  the breakpoint are underlined; since  $\mathcal{A}(G, H) = 1$  and  $\mathcal{T}(G, H) = 1$  and  $\mathcal{G}(G) = 10$ ,  $\mathcal{D}_{BP}(G, H) = 8.5$ .

### 4. QUASI MAXIMUM PERFECT AND COVERING SETS

Since polyploids are perfect genomes where every gene family has the same size, our strategy to aliquote a genome  $G$  is to find a perfect subset of  $G$  and then transform it into a perfect genome. Because manipulating  $G$  directly is challenging, we will manipulate its unordered counterpart  $\tilde{G}$  instead.

There are many possible perfect subsets of  $\tilde{G}$  and we need to choose the one which will minimize the distance. From Equation 1, we can see that the distance is minimized by maximizing the number of adjacencies and telomeres present in both genomes. Thus, we define the *weight* of a perfect set  $P$  in relation to a genome  $\tilde{G}$  to be:

$$\mathcal{W}(P) = \mathcal{A}(\tilde{G}, P) + \frac{\mathcal{T}(\tilde{G}, P)}{2} \quad (2)$$

We define a *maximum perfect set*  $P$  to be a perfect subset of  $\tilde{G}$  with maximum weight.

**Algorithm 1:** Replicate

---

**Input:**  $P$  quasi maximum perfect and covering set with respect to a genome  $\tilde{G}$  and  $p$ , the size of the gene families in  $\tilde{G}$ .  
**Output:** A polyploid  $\tilde{H}$  with all gene families of size  $p$ .

- 1  $\tilde{H} = P$
- 2 **foreach**  $\alpha \in P$  **do**
- 3   **if** the multiplicity of  $\alpha$  is not  $p$  **then**
- 4     add  $\alpha$  to  $\tilde{H}$  until the multiplicity is  $p$
- 5   **end**
- 6 **end**
- 7 **return**  $\tilde{H}$

---

While it is impossible to compute the distance using an unordered genome, it is easy to see that, once transformed into a genome and reordered, a maximum perfect subset of  $\tilde{G}$  will minimize the distance. However, not all maximum perfect sets can be transformed into an ordered genome that can be compared with  $G$ . In addition to being ordered, to be compared two genomes must have the same set of extremities and it is possible to construct a maximum perfect set that is missing extremities from  $\tilde{G}$ . Thus, we introduce the notion of a *covering set*. A set  $C$  covers an unordered genome  $\tilde{G}$  if and only if it contains at least one copy of every extremity in  $\tilde{G}$ .

Unfortunately, it is impossible to always find a maximum perfect and covering subset of  $\tilde{G}$ ; for some unordered genomes no such subset exists. To solve the problem, we simply remove the restriction that it must be a subset, i.e., the maximum perfect and covering set may have elements not in  $\tilde{G}$ . However, without the restriction that the perfect set has to be a subset of  $\tilde{G}$ , the idea of a maximum perfect set does not make any sense; we could add an infinite number of adjacencies and telomeres to the set if they do not have to come from the genome. Thus, we introduce a *quasi maximum perfect set*: a perfect set with a maximum perfect set as a subset but with additional elements not from  $\tilde{G}$ . All maximum perfect sets are also quasi maximum perfect sets, although the reverse is not true.

It is easy to transform a quasi maximum perfect and covering set with respect to a genome  $\tilde{G}$  into a polyploid that, if ordered, could be compared with  $G$ : we simply increase the number of copies of each element. Algorithm 1 does exactly this; its correctness follows from the following lemma:

**Lemma 1.** *Given a genome  $\tilde{G}$  with all gene families of size  $p$  and  $P$ , a quasi maximum perfect and covering set with respect to  $\tilde{G}$ ,  $\tilde{H} = \text{REPLICATE}(P, p)$  is a polyploid such that  $\mathcal{E}(\tilde{H}) = \mathcal{E}(\tilde{G})$ .*

**Proof.** There are three properties of  $\tilde{H}$  that we need to prove. First, we need to prove that  $\tilde{H}$  is in fact an unordered genome. Second, that it is a polyploid. Third, that  $\mathcal{E}(\tilde{H}) = \mathcal{E}(\tilde{G})$ .

The fact that  $\tilde{H}$  is unordered is trivial.  $\tilde{H}$  is a genome follows from the fact that  $P \subseteq \tilde{H}$  and  $P$  covers  $\tilde{G}$  and Algorithm 1 ensures an equal number of heads and tails.  $\mathcal{E}(\tilde{H}) = \mathcal{E}(\tilde{G})$  also follows from this fact. Hence,  $\tilde{H}$  is an unordered genome and  $\mathcal{E}(\tilde{H}) = \mathcal{E}(\tilde{G})$ .

By definition,  $\tilde{H}$  is a polyploid if and only if it is perfect with all gene families of size  $p$ . From Algorithm 1, it is clear that all gene families of  $\tilde{H}$  are of size  $p$ . Since  $P$  is perfect and  $P \subseteq \tilde{H}$  and, furthermore, the only difference between  $P$  and  $\tilde{H}$  is the multiplicity of the elements, which does not affect the property of perfection,  $\tilde{H}$  must also be perfect. ■

## 5. REORDER GENES

How a genome is reordered can affect the distance. However, with a quasi maximum perfect and covering set with respect to  $\tilde{G}$  as a base, it is easy to reorder the resulting polyploid such that its distance to  $G$  is minimal. Algorithm 2 reorders the genome, and the following lemma proves its correctness:

**Lemma 2.** *Given a genome  $G$  with all gene families of size  $p$ , let  $P$  be a quasi maximum perfect and covering set with respect to  $\tilde{G}$  and let  $\tilde{H} = \text{REPLICATE}(P, p)$ , then  $\mathcal{D}_{BP}(G, H)$  where  $H = \text{REORDER}(\tilde{H}, G)$  is minimal.*

---

**Algorithm 2:** Reorder

---

**Input:** A polyploid  $\tilde{H}$  with all gene families of size  $p$  and a genome  $G$  with all gene families of size  $p$  where  $\mathcal{E}(\tilde{G}) = \mathcal{E}(\tilde{H})$ .

**Output:** A polyploid  $H$  with all gene families of size  $p$  such that  $\mathcal{D}_{BP}(G, H)$  is minimal among all ordered genomes of  $\tilde{H}$ .

```

1   $H \leftarrow \emptyset$ 
2  foreach element  $\alpha \in G$  do
3    if  $\tilde{\alpha} \in \tilde{H}$  then
4      add  $\alpha$  to  $H$ 
5    end
6  end
7   $E \leftarrow \mathcal{E}(G) \setminus \mathcal{E}(H)$ 
8  foreach element  $\tilde{\alpha} \in \tilde{H}$  do
9    if  $\tilde{\alpha}$  is an adjacency then
10   let  $\{x, y\} = \tilde{\alpha}$  where  $x, y$  are extremities
11   while there exists an  $i$  such that  $x_i \in E$  and a  $j$  such that  $y_j \in E$  do
12     add  $\{x_i, y_j\}$  to  $H$ 
13     remove  $x_i$  from  $E$ 
14     remove  $y_j$  from  $E$ 
15   end
16   else
17     let  $\{x\} = \tilde{\alpha}$  where  $x$  is an extremity
18     while there exists an  $i$  such that  $x_i \in E$  do
19       add  $\{x_i\}$  to  $H$ 
20       remove  $x_i$  from  $E$ 
21     end
22   end
23 end
24 return  $H$ 

```

---

**Proof.**  $\mathcal{D}_{BP}(G, H)$  is minimal if there does not exist an other genome  $H'$  such that  $\mathcal{D}_{BP}(G, H') < \mathcal{D}_{BP}(G, H)$ . Assume towards contradiction that such an  $H'$  does exist.

Let  $E$  be the set of elements that are identical between  $H$  and  $G$  and let  $E'$  be the set of elements that are identical between  $H'$  and  $G$ . We observe that because  $H$  and  $H'$  are polyploids  $\tilde{E}$  and  $\tilde{E}'$  must be perfect.

From Equation 1, only adjacencies and telomeres that are the same in both genomes reduce the distance. Thus, since  $\mathcal{D}_{BP}(G, H') < \mathcal{D}_{BP}(G, H)$ ,  $\mathcal{W}(\tilde{E}') > \mathcal{W}(\tilde{E})$ .

By definition, a quasi maximum perfect set must have a maximum perfect set as a subset; let  $M \subseteq P$  be such a maximum perfect set of  $\tilde{G}$ . Since  $E$  is a subset of  $H$ , it follows that  $\tilde{E}$  must be a subset of  $\tilde{H}$ . Thus,  $M$  and  $\tilde{E}$  are subsets of  $\tilde{H}$  and, in fact, we will prove that  $M = \tilde{E}$  by assuming towards contradiction that  $M \neq \tilde{E}$ .

$M \neq \tilde{E}$  if and only if the multiplicity of all elements in  $M$  is not equal to the multiplicity of all elements in  $\tilde{E}$ . Let  $\alpha$  be an element that has different multiplicities in both  $M$  and  $\tilde{E}$ . Since  $M$  is a maximum perfect set of  $\tilde{G}$  all of its elements can have a multiplicity of at most  $p$  and, because  $M$  is maximum, the multiplicity of  $\alpha$  in  $M$  is equal to the multiplicity of  $\alpha$  in  $\tilde{G}$ . Since REPLICATE( $P, p$ ) creates  $\tilde{H}$  by setting all the elements of  $P$ , which includes all of the elements of  $M$ , to multiplicity  $p$ , the multiplicity of  $\alpha$  in  $\tilde{H}$  is greater than or equal to that of  $M$ . It follows that the multiplicity of  $\alpha$  in  $M$  is the multiplicity of  $\alpha$  in  $\tilde{G} \cap \tilde{H}$ . From Lines 2–6 in Algorithm 2 REORDER( $\tilde{H}, G$ ) we can conclude that  $\tilde{G} \cap \tilde{H} = (G \cap H)$ . By definition,  $E = G \cap H$ , hence,  $(G \cap H) = \tilde{E}$ . Therefore, the multiplicities of  $\alpha$  in  $M$  and in  $\tilde{E}$  are equal; a contradiction. Thus,  $M = \tilde{E}$ .

$\mathcal{W}(\tilde{E}') > \mathcal{W}(M)$  follows from the facts that  $M = \tilde{E}$  and  $\mathcal{W}(\tilde{E}') > \mathcal{W}(\tilde{E})$ . But  $M$  is a maximum perfect set; a contradiction. Hence,  $H$  must be minimal. ■

## 6. FULLY MODIFIED CLIQUE GRAPHS

We have reduced the problem of aliquoting a genome to that of finding a maximum perfect and covering set in the breakpoint model. To find the maximum perfect and covering set, we construct a graph from the genome that highlights its important information.

**Definition 4.** A fully modified weighted clique graph of a genome  $\tilde{G}$ , denoted  $CG(\tilde{G})$ , is defined as follows:

- For each gene family  $x$  in  $\tilde{G}$ ,  $CG(\tilde{G})$ , has four vertices, one labeled  $\vec{x}$ , one labeled  $\overleftarrow{x}$  and two without labels, called null vertices, of which one is connected to the vertex labeled with  $\vec{x}$  and the other to the vertex labeled with  $\overleftarrow{x}$ ;
- For each adjacency  $\{x, y\}$  in  $\tilde{G}$  where  $x \neq y$ , there is an edge between the vertices labeled  $x$  and  $y$  with a weight equal to the multiplicity of  $\{x, y\}$  in  $\tilde{G}$ ;
- For each telomere  $\{x\}$  in  $\tilde{G}$ , the edge between the vertex labeled  $x$  and its adjacent null vertex has a weight equal to half of the multiplicity of  $\{x\}$  in  $\tilde{G}$ ;
- All other edges in  $CG(\tilde{G})$  have weights of 0.

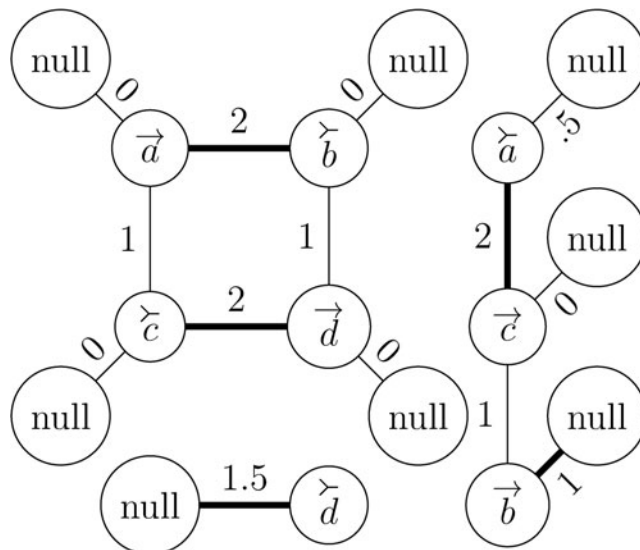
For example, the modified weighted clique graph of the genome  $\tilde{G} = \{\{\vec{a}\}, \{\vec{a}, \vec{b}\}, \{\vec{b}, \vec{c}\}, \{\vec{c}, \vec{a}\}, \{\vec{a}, \vec{c}\}, \{\vec{c}, \vec{d}\}, \{\vec{d}\}, \{\vec{d}, \vec{c}\}, \{\vec{c}, \vec{a}\}, \{\vec{a}, \vec{b}\}, \{\vec{b}\}, \{\vec{b}\}, \{\vec{b}, \vec{d}\}, \{\vec{d}\}\}$  is depicted in Figure 1.

In the fully modified weighted clique graph, edges correspond to adjacencies or telomeres. Thus, the problem is to choose edges from the fully modified weighted clique graph such that the corresponding adjacencies and telomeres form a maximum perfect and covering set. From Definition 4, it is clear that the adjacencies or telomeres corresponding to any two edges in the fully modified weighted clique graph are perfect if and only if the edges share no vertex in common or correspond to the same edge (since there are no loops, there are no edges that correspond to adjacencies of the form  $\{x, x\}$  for some extremity  $x$ ). Finding a set of edges that share no vertices in common is a well-studied computational problem known as the *maximum matching problem*, although, since our edges are weighted, in this case we are actually more interested in the *maximum weight matching problem* Lovász and Plummer (2009).

Given a matching, Algorithm 3 constructs its corresponding perfect set. However, it is the quality of the matching that determines the quality of the perfect set. The fully modified weighted clique graph was defined such that the weight of the maximum weight matching is equal to the weight of the perfect set defined by Algorithm 3. Thus, given a maximum weight matching, Algorithm 3 will construct a maximum perfect set:

**Lemma 3.** Given a genome  $\tilde{G}$ , let  $M$  be a matching of  $CG(\tilde{G})$  and let  $P = \text{PERFECTSET}(\tilde{G}, M)$ . If  $M$  is a maximum weighted matching then  $P$  must be a quasi maximum perfect set.

**Proof.** Assume towards contradiction that  $M$  is a maximum weight matching but  $P$  is not a quasi maximum perfect set. Then there must exist a perfect set  $P'$  such that  $\mathcal{W}(P) < \mathcal{W}(P')$ . Let  $M'$  be the matching that correspond to  $P'$ . Since the weight of the perfect set and its corresponding matching are the same,  $\mathcal{W}(M) < \mathcal{W}(M')$ . But  $M$  is a maximum weight matching, a contradiction. ■



**FIG. 1.** An example of a fully modified weighted clique graph. The five bold edges represent all the edges that belong to the maximum weight matching of this graph.

**Algorithm 3:** PerfectSet

---

**Input:** A genome  $\tilde{G}$  and a matching  $M$ .  
**Output:** A perfect set  $P$ .

```

1   $P \leftarrow \emptyset$ 
2  foreach edge  $e \in M$  do
3    let  $e = \{u, v\}$  for vertices  $u$  and  $v$ 
4    if  $v$  is a null vertex then
5      add the telomere that corresponds to  $u$  to  $P$ 
6    else if  $u$  is a null vertex then
7      add the telomere that corresponds to  $v$  to  $P$ 
8    else
9      add the adjacencies that corresponds to  $u$  and  $v$  to  $P$ 
10   end
11 end
12 return  $P$ 

```

---

While in many cases the maximum weight matching will correspond to a quasi maximum perfect set that also covers the genome, simply finding a maximum weight matching of the fully modified weighted clique graph does not guarantee that the perfect set will be covering. Fortunately, any matching can be easily modified by Algorithm 4 so that this is the case:

**Lemma 4.**  $M = \text{EXTENDED MAXIMUM WEIGHT MATCHING}(\mathcal{CG}(\tilde{G}))$ , where  $\tilde{G}$  is an unordered genome, is a maximum weight matching such that  $P = \text{PERFECT SET}(\tilde{G}, M)$  covers  $\tilde{G}$ .

**Proof.** Since every non-null vertex corresponds to an extremity in a genome, the perfect set will only be covering if every non-null vertex is incident with an edge in the matching from which the perfect set was created.

Algorithm 4 first finds any maximum weight matching  $M'$ . If every non-null vertex is incident with an edge in  $M'$  then the algorithm does not modify the matching and the result is a maximum weight matching whose corresponding perfect set covers  $\tilde{G}$ .

Each non-null vertex has a null vertex to which only it is adjacent. For non-null vertices not incident in an edge in the maximum weight matching, Algorithm 4 adds the edge between this vertex and its null counterpart to the matching. Clearly, the resulting perfect set will cover  $\tilde{G}$ . Thus, we must ensure that the resulting set of edges is still a matching of maximum weight.

The resulting set of edges must be a matching since the non-null vertex is not incident with any edges in the matching and the null vertex is only adjacent with the non-null vertex. It must be a maximum weight matching since all weights in a fully modified weighted clique graph are non-negative.

Therefore,  $M = \text{EXTENDED MAXIMUM WEIGHT MATCHING}(\mathcal{CG}(\tilde{G}))$  is a maximum weight matching such that  $P = \text{PERFECT SET}(\tilde{G}, M)$  covers  $\tilde{G}$ . ■

**Algorithm 4:** ExtendedMaximumWeightMatching

---

**Input:** A fully modified weighted clique graph  $\mathcal{CG}(\tilde{G})$  of a genome  $\tilde{G}$ .  
**Output:** A maximum weight matching  $M$  of  $\mathcal{CG}(\tilde{G})$  where every non-null vertex in  $\mathcal{CG}(\tilde{G})$  is incident with an edge in  $M$ .

```

1   $M \leftarrow \text{MAXIMUM WEIGHT MATCHING}(\mathcal{CG}(\tilde{G}))$ 
2  foreach non-null vertex  $v$  in  $\mathcal{CG}(\tilde{G})$  do
3    if  $v$  is not incident with an edge in  $M$  then
4      let  $u$  be the null vertex adjacent to  $v$ 
5      add  $\{u, v\}$  to  $M$ 
6    end
7  end
8  return  $M$ 

```

---

**Algorithm 5:** BreakpointAliquoting

---

**Input:** A genome  $G$  with all gene families of size  $p$ .  
**Output:** A polyploid  $H$  with all gene families of size  $p$ .

- 1  $M \leftarrow \text{EXTENDED\_MAXIMUM\_WEIGHT\_MATCHING}(\mathcal{CG}(\tilde{G}))$
- 2  $P \leftarrow \text{PERFECT\_SET}(\tilde{G}, M)$
- 3  $\tilde{H} \leftarrow \text{REPLICATE}(P, p)$
- 4  $H \leftarrow (\tilde{H}, G)$
- 5 **return**  $H$

---

**7. ALGORITHM**

From Lemmas 3 and 4, we can conclude the following theorem:

**Theorem 1.** *Let  $M = \text{EXTENDED\_MAXIMUM\_WEIGHT\_MATCHING}(\mathcal{CG}(\tilde{G}))$  and  $P = \text{PERFECT\_SET}(\tilde{G}, M)$ , where  $\tilde{G}$  is an unordered genome.  $P$  is a quasi maximum perfect and covering set of  $\tilde{G}$ .*

Algorithm 5 is the complete breakpoint aliquoting algorithm bringing together all the algorithms discussed in the previous sections. It follows from Lemma 2 and Theorem 1 that Algorithm 5 produces the optimal breakpoint aliquoting.

Most of the algorithms discussed in this article run in linear time, however, the best known algorithm for computing the maximum weight matching runs in  $O(ev \log v)$  time (Lovász and Plummer, 2009) where  $e$  is the number of edges and  $v$  is the number of vertices. If we have a genome with  $n$  gene families of size  $p$ , the fully modified weighted clique graph will have  $4n$  vertices and  $(2p - 2)n$  edges. Thus, our algorithm runs in  $O(pn^2 \log n)$  time.

**8. EXAMPLE**

Suppose the genome  $G = \{\{\tilde{a}_1\}, \{\tilde{a}_1, \tilde{b}_1\}, \{\tilde{b}_1, \tilde{c}_1\}, \{\tilde{c}_1, \tilde{a}_2\}, \{\tilde{a}_2, \tilde{c}_2\}, \{\tilde{c}_2, \tilde{a}_1\}, \{\tilde{d}_1\}, \{\tilde{d}_2\}, \{\tilde{d}_2, \tilde{c}_3\}, \{\tilde{c}_3, \tilde{a}_3\}, \{\tilde{a}_3, \tilde{b}_2\}, \{\tilde{b}_2\}, \{\tilde{b}_3\}, \{\tilde{b}_3, \tilde{d}_3\}, \{\tilde{d}_3\}\}$  is given as input to Algorithm 5. All the gene families of  $G$  contain three genes, hence, we hope to reconstruct a *hexaploid* ancestor, an ancestor with three different types chromosomes for a total of six chromosomes, of  $G$  using the algorithm. Its corresponding unordered genome is  $\tilde{G} = \{\{\tilde{a}\}, \{\tilde{a}, \tilde{b}\}, \{\tilde{b}, \tilde{c}\}, \{\tilde{c}, \tilde{a}\}, \{\tilde{a}, \tilde{c}\}, \{\tilde{c}, \tilde{d}\}, \{\tilde{d}\}, \{\tilde{d}\}, \{\tilde{d}, \tilde{c}\}, \{\tilde{c}, \tilde{a}\}, \{\tilde{a}, \tilde{b}\}, \{\tilde{b}\}, \{\tilde{b}\}, \{\tilde{b}, \tilde{d}\}, \{\tilde{d}\}\}$  and its fully modified weighted clique graph is depicted in Figure 1.

From Figure 1, we can see that the maximum weight matching of  $\mathcal{CG}(\tilde{G})$  consists of the set of edges  $M = \{\{\tilde{a}, \tilde{b}\}, \{\tilde{c}, \tilde{d}\}, \{\text{null}, \tilde{d}\}, \{\tilde{a}, \tilde{c}\}, \{\text{null}, \tilde{b}\}\}$ . As is often the case, in this particular example, all the extremities are present so extending the matching with Algorithm 4 adds nothing to  $M$ .

**Remark 1.** *As an interesting aside, we observe that at this point it is possible to calculate the distance between  $G$  and its most parsimonious polyploid ancestor. While the perfect set is not computed until the next step of the algorithm, the weight of the perfect set is known now. From Lemma 3, the sum of weights of the edges of  $\mathcal{CG}(\tilde{G})$  is equivalent to the weight of the perfect set. Since the weight of the perfect set is based on the breakpoint distance equation, we can insert the weight of the matching perfect set into the equation by subtracting it from  $\mathcal{G}(\tilde{G})$ . The result is the breakpoint distance between  $G$  and its most parsimonious polyploid ancestor.*

We can calculate the final breakpoint distance given the weight of  $M$ . Since  $\mathcal{W}(M) = 8.5$  and  $\mathcal{G}(G) = 12$ , the breakpoint distance between  $G$  and the hexaploid we will construct will be 3.5.

Applying Algorithm 3 to  $M$  produces the perfect set  $P = \{\{\tilde{a}, \tilde{b}\}, \{\tilde{c}, \tilde{d}\}, \{\tilde{d}\}, \{\tilde{a}, \tilde{c}\}, \{\tilde{b}\}\}$ . Algorithm 1 then increases the sizes of the gene families in  $P$  to produce the unordered hexaploid  $\tilde{H} = \{\{\tilde{b}\}\{\tilde{a}, \tilde{b}\}, \{\tilde{a}, \tilde{c}\}, \{\tilde{c}, \tilde{d}\}, \{\tilde{d}\}, \{\tilde{b}\}\{\tilde{a}, \tilde{b}\}, \{\tilde{a}, \tilde{c}\}, \{\tilde{c}, \tilde{d}\}, \{\tilde{d}\}, \{\tilde{b}\}, \{\tilde{a}, \tilde{b}\}, \{\tilde{a}, \tilde{c}\}, \{\tilde{c}, \tilde{d}\}, \{\tilde{d}\}\}$ .

Finally, Algorithm 2 is applied to  $\tilde{H}$ .  $\{\tilde{d}\}$  in  $\tilde{H}$  has three ordered equivalents in  $G$ :  $\{\tilde{d}_1\}, \{\tilde{d}_2\}, \{\tilde{d}_3\}$ . The rest of the adjacencies in  $\tilde{H}$  each have two ordered equivalents in  $G$ , for example,  $\{\tilde{c}, \tilde{d}\}$  is equivalent to



$\{\vec{c}_2, \vec{d}_1\}$  and  $\{\vec{d}_2, \vec{c}_3\}$ . Since all the gene families are of size three and for every adjacency two of the three ordered adjacencies have equivalents in  $G$ , it is trivial to calculate the third ordered adjacency, in the case of  $\{\vec{c}, \vec{d}\}$  it is  $\{\vec{c}_1, \vec{d}_3\}$ . However, even in more complex examples, the remaining adjacencies are just as easy to calculate as the pairing of the unused ordered extremities is arbitrary. The resulting ordered hexaploid is  $H = \{\{\vec{b}_1\}, \{\vec{a}_1, \vec{b}_1\}, \{\vec{a}_1, \vec{c}_1\}, \{\vec{c}_1, \vec{d}_3\}, \{\vec{d}_3\}, \{\vec{b}_2\}, \{\vec{a}_3, \vec{b}_2\}, \{\vec{c}_3, \vec{d}_3\}, \{\vec{d}_2, \vec{c}_3\}, \{\vec{d}_2\}, \{\vec{b}_3\}, \{\vec{a}_2, \vec{b}_3\}, \{\vec{a}_2, \vec{c}_2\}, \{\vec{c}_2, \vec{d}_1\}, \{\vec{d}_1\}\}$ .

Since both  $H$  and  $G$  are ordered we can compare them. The only breakpoints are in  $H$  are  $\{\vec{b}_1\}, \{\vec{a}_1, \vec{c}_1\}, \{\vec{c}_1, \vec{d}_3\}$  and  $\{\vec{a}_2, \vec{b}_3\}$ ; since  $\mathcal{A}(G, H) = 6$  and  $\mathcal{T}(G, H) = 5$  and  $\mathcal{G}(G) = 12$ , from Equation 2,  $\mathcal{D}_{BP} = 3.5$ . Observe that this matches perfectly with our earlier calculation.

### 9. DOUBLE CUT AND JOIN DISTANCE

*Double cut and join* distance and breakpoint distance are closely related. BP distance counts the number of differences between two genomes. For the most part, this is the same as DCJ distance; it generally takes one DCJ operation to correct one difference between two genomes. However, sometimes a DCJ operation corrects two differences instead of one, and the DCJ distance must take this into account; this is the only difference between the two distances and it is not very significant as this occurs infrequently with most data sets encountered in practice.

**Definition 5** (Bergeron et al., 2006). *The double cut and join (DCJ) operation acts on two adjacencies or telomeres  $u$  and  $v$  of a genome in one of the following three ways:*

- *If both  $u = \{w, x\}$  and  $v = \{y, z\}$  are adjacencies, these are replaced by the two adjacencies  $\{w, y\}$  and  $\{x, z\}$  or by the two adjacencies  $\{w, z\}$  and  $\{x, y\}$ .*
- *If  $u = \{w, x\}$  is an adjacency and  $v = \{y\}$  is a telomere, these are replaced by  $\{w, y\}$  and  $\{x\}$  or by  $\{x, y\}$  and  $\{w\}$ .*
- *If both  $u = \{w\}$  and  $v = \{x\}$  are telomeres, these are replaced by  $\{w, x\}$ .*

*In addition, as an inverse of the last case, a single adjacency  $\{w, x\}$  can be replaced by two telomeres  $\{w\}$  and  $\{x\}$ .*

DCJ distance is easily computed with the aid of the following graph:

**Definition 6** (Bergeron et al., 2006). *An adjacency graph (depicted in Fig. 2) is a bipartite graph with two sets of vertices labeled with the elements of the two genomes that are to be compared respectively. There is an edge connecting two vertices for each extremity in common between their corresponding adjacencies.*

Recall that to be compared using a distance algorithm like DCJ, both genomes must have the same genes, hence the same extremities, and, in the case of genomes with duplicated genes, must be ordered.

From the adjacency graph, the DCJ distance can be computed:

$$\mathcal{D}_{DCJ}(G, H) = |\mathcal{G}(G)| - c - \frac{i}{2} \tag{3}$$

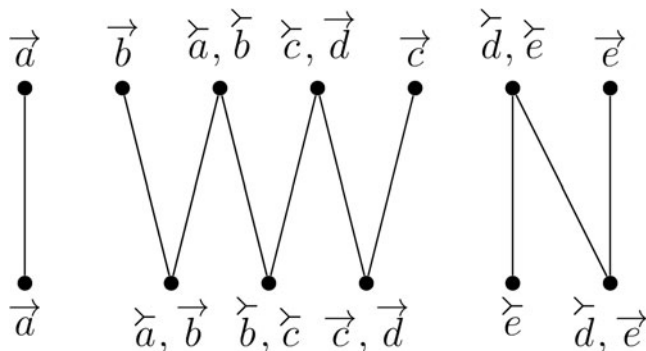


FIG. 2. An example of an adjacency graph.

where  $c$  is the number of cycles and  $i$  is the number of paths with an odd number of edges in the adjacency graph.

The following theorem gives the relationship between the DCJ distance and BP distance:

**Theorem 2.** *Given two genomes  $G$  and  $H$ :*

$$\mathcal{D}_{DCJ}(G, H) \leq \mathcal{D}_{BP}(G, H) \leq 2 \cdot \mathcal{D}_{DCJ}(G, H) \tag{4}$$

**Proof.** First, for  $\mathcal{D}_{DCJ}(G, H) \leq \mathcal{D}_{BP}(G, H)$ , observe that every adjacency shared between the two genomes causes a cycle of size 2 in the adjacency graph. It follows that  $|\mathcal{A}(G, H)|$  is equal to the number of cycles of size 2 in the adjacency graph and, hence, less than or equal to the number of cycles in the adjacency graph. By similar reasoning,  $|\mathcal{T}(G, H)|$  is equal to the number of paths of size 1 in the adjacency graph and, hence, less than or equal to the number of paths in the adjacency graph. It follows that  $|\mathcal{A}(G, H)| + \frac{|\mathcal{T}(G, H)|}{2} \leq c + \frac{i}{2}$  where  $c$  is the number of cycles and  $i$  the number of odd paths in the adjacency graph. Since this factor decreases the distance, it follows that  $\mathcal{D}_{DCJ}(G, H) \leq \mathcal{D}_{BP}(G, H)$ .

Second, for  $\mathcal{D}_{BP}(G, H) \leq 2 \cdot \mathcal{D}_{DCJ}(G, H)$ , we will assume the most extreme case, where  $|\mathcal{A}(G, H)| + \frac{|\mathcal{T}(G, H)|}{2} = 0$  but the number of cycles and paths in the adjacency graph is otherwise maximal. This produces the worst possible BP distance,  $|\mathcal{G}(G, H)|$ , but it is otherwise a maximal DCJ distance.

To determine the maximum number of cycles and paths in the adjacency graph, observe that every edge in the adjacency graph corresponds to one extremity shared between the genomes. Since  $|\mathcal{A}(G, H)| + \frac{|\mathcal{T}(G, H)|}{2} = 0$  it follows that there are no cycles of size 2 and odd paths of size 1, so the smallest cycles and odd paths are of size 4 (since all cycles are even) and 3 respectively. Thus, the maximum number of cycles and odd paths is  $\frac{|\mathcal{G}(G, H)|}{2}$ . Thus, in the worst case scenario, BP distance is equal to twice the DCJ distance. ■

Because BP distance can be equal to the DCJ distance, it is possible that the genome that results from Algorithm 5 could represent the optimal aliquoting for both BP and DCJ. Because  $\mathcal{D}_{BP}(G, H) \leq 2 \cdot \mathcal{D}_{DCJ}(G, H)$ , the distance between the original genome and its optimal BP aliquoting is no more than twice the distance between the original genome and its optimal DCJ distance. Thus, we can deduce the following:

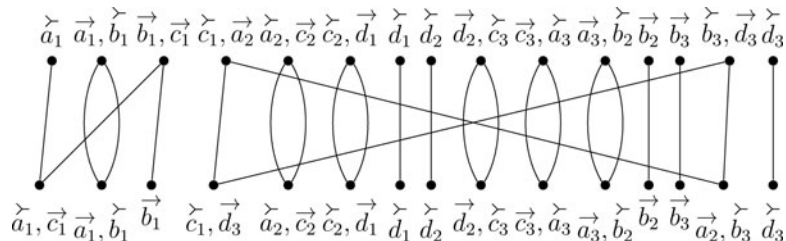
**Theorem 3.** *Algorithm 5 is a 2-approximation for the genome aliquoting problem using DCJ distance.*

We can see that Theorem 3 holds true in the case of the example discussed in Section 7. Figure 3 depicts the adjacency graph between the genomes  $G$  and  $H$  from Section 7. As we can see,  $c=7$ ,  $i=6$  and  $\mathcal{G}(G) = 12$ ; hence,  $\mathcal{D}_{DCJ}(G, H) = 2$ . While  $\mathcal{D}_{BP}(G, H) = 3.5$ , close to the upper bound of twice  $\mathcal{D}_{DCJ}(G, H)$ , the resulting double cut and join distance is very good and, in fact, is almost certainly optimal.

### 10. CONCLUSION

Since the input to our algorithm is of size  $pn$  where  $p$  is the size of the gene families and  $n$  is the number of gene families, our algorithm runs in sub-cubic time. Thus, there exists a polynomial time algorithm that solves the genome aliquoting problem for BP distance. However, our algorithm is not without limitations.

**FIG. 3.** An adjacency graph between an order genome with all gene families of size three (top) and one possible hexaploid ancestor (bottom) where the distance between the two genomes is minimal.



The output of our algorithm consists of a mixture of linear and circular chromosomes with the only restriction being that adjacencies of the form  $\{x, x\}$  for some extremity  $x$  are not permitted. In many cases, it is desirable to restrict the output to something more similar to the input. For example, if the input consists of linear chromosomes then the output should consist of linear chromosomes. We believe it is possible to modify the output of our algorithm to provide the desired output in most cases without changing the distance. Some of the genome halving algorithms do precisely this (El-Mabrouk and Sankoff, 2003; Alekseyev and Pevzner, 2007). On the other hand, sometimes circular chromosomes with adjacencies of the form  $\{x, x\}$  are desirable in the output, in which case it might be possible, in some cases, to get a better result than produced by our algorithm.

It remains an open problem whether or not a polynomial time algorithm for the genome aliquoting problem with DCJ distance exists. However, as we have shown, a good approximation algorithm does exist. We remain optimistic that a polynomial time algorithm for DCJ distance can be found in the future.

### ACKNOWLEDGMENTS

Research was supported by a Discovery grant from the Natural Sciences and Engineering Research Council of Canada. David Sankoff holds the Canada Research Chair in Mathematical Genomics.

### DISCLOSURE STATEMENT

No competing financial interests exist.

### REFERENCES

- Alekseyev, M.A., and Pevzner, P.A. 2007. Whole genome duplications and contracted breakpoint graphs. *SIAM J. Comput.*, 36, 1748–1763.
- Bergeron, A., Mixtacki, J., and Stoye, J. 2006. A unifying view of genome rearrangements. *Lect. Notes Comput. Sci.* 4175, 163–173.
- El-Mabrouk, N., and Sankoff, D. 2003. The reconstruction of doubled genomes. *SIAM J. Comput.* 32, 754–792.
- Feijão, P., and Meidanis, J. 2009. SCJ: a novel rearrangement operation for which sorting, genome median and genome halving problems are easy. *Lect. Notes Comput. Sci.* 5724, 85–96.
- Lovász, L., and Plummer, M.D. 2009. *Matching Theory*. AMS Chelsea Publishing, Providence.
- Mixtacki, J. 2008. Genome halving under DCJ revisited. *Lect. Notes Comput. Sci.* 5092, 276–286.
- Sankoff, D., and Blanchette, M. 1997. The median problem for breakpoints in comparative genomics. *Proc. COCOON '97* 251–264.
- Tannier, E., Zheng, C., and Sankoff, D. 2009. Multichromosomal median and halving problems under different genomic distances. *Bioinformatics* 10, 120.
- Warren, R., and Sankoff, D. 2009a. Genome aliquoting with double cut and join. *BMC Bioinform.* 10, S2.
- Warren, R., and Sankoff, D. 2009b. Genome halving with double cut and join. *J. Bioinform. Comput. Biol.* 7, 357–371.

Address correspondence to:

Dr. David Sankoff  
Department of Mathematics and Statistics  
University of Ottawa  
585 King Edward Avenue  
Ottawa, Canada

E-mail: sankoff@uottawa.ca

