



HAL
open science

A General Framework for Genome Rearrangement with Biological Constraints

Pijus Simonaitis, Annie Chateau, Krister M. Swenson

► **To cite this version:**

Pijus Simonaitis, Annie Chateau, Krister M. Swenson. A General Framework for Genome Rearrangement with Biological Constraints. RECOMB-CG: Comparative Genomics, Oct 2018, Magog-Orford, QC, Canada. pp.49-71, 10.1007/978-3-030-00834-5_3 . lirmm-02067487

HAL Id: lirmm-02067487

<https://hal-lirmm.ccsd.cnrs.fr/lirmm-02067487>

Submitted on 17 Sep 2019

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - ShareAlike 4.0 International License

RESEARCH

Open Access



A general framework for genome rearrangement with biological constraints

Pijus Simonaitis¹, Annie Chateau^{1,2} and Krister M. Swenson^{1,2*} 

Abstract

This paper generalizes previous studies on genome rearrangement under biological constraints, using double cut and join (DCJ). We propose a model for weighted DCJ, along with a family of optimization problems called φ -MCPS (MINIMUM COST PARSIMONIOUS SCENARIO), that are based on labeled graphs. We show how to compute solutions to general instances of φ -MCPS, given an algorithm to compute φ -MCPS on a circular genome with exactly one occurrence of each gene. These general instances can have an arbitrary number of circular and linear chromosomes, and arbitrary gene content. The practicality of the framework is displayed by presenting polynomial-time algorithms that generalize the results of Bulteau, Fertin, and Tannier on the SORTING BY WDCJS AND INDELS IN INTERGENES problem, and that generalize previous results on the MINIMUM LOCAL PARSIMONIOUS SCENARIO problem.

Keywords: Double cut and join (DCJ), Weighted genome rearrangement, Breakpoint graph, Graph edit distance, Edge switch

Introduction

Context

The practical study of genome rearrangement scenarios has been limited by a lack of mathematical models capable of incorporating biological constraints, since foundational models focused on minimum length scenarios transforming one genome into another. In the modern age, where the collection of fully assembled and annotated genomes is ever-increasing, there is the need for the development of more elaborate mathematical models that consider the data from multiple biological experiments.

One way to incorporate biological information into the inference of evolutionary scenarios is to consider models that weight rearrangements according to their likelihood of occurring; a breakpoint may be more likely to occur in some intergenic regions than others. To this end, the study of length-weighted reversals was started in the late nineties by Blanchette et al. [1]. Baudet et al. present a summary of work done in this area, along with

work on reversals centered around the origin of replication [2]. Recently, Tannier has published a series of papers focused on weighting intergenic regions by their length in nucleotides. In [3], Biller et al. pointed out that, according to the Nadeau–Taylor model of uniform random breakage [4, 5], a breakpoint is more likely to occur in a longer intergenic region. Subsequent papers by Fertin et al. [6], and Bulteau et al. [7] present algorithmic results for models that take into account the length of intergenic regions. Using Hi-C data [8], Veron et al. along with our own study, have pointed out the importance of weighting pairs of breakpoints according to how close they tend to be in physical space [9, 10]. In order to use this physical constraint, we partitioned intergenic regions into co-localized areas, and developed algorithms for computing distances that minimize the number of rearrangements that operate on breakpoints between different areas [11, 12].

Much of this work is based on the mathematically clean model for genome rearrangement called *Double Cut and Join*, or *DCJ* [13, 14]. Genomes are partitioned into n orthologous syntenic blocks that we will simply call *genes*. Each gene is represented by two extremities, and each chromosome is represented by an ordering of these extremities. Those extremities that are adjacent in

*Correspondence: swenson@lirmm.fr

¹ CNRS, LIRMM, Université Montpellier, 161 Rue Ada, 34392 Montpellier, France

Full list of author information is available at the end of the article



this ordering are paired, and transformations of these pairs occur by swapping extremities of two pairs. DCJ can naturally be interpreted as a graph edit model with the use of the *breakpoint graph*, where there is an edge between gene extremities a and b for each adjacent pair. A DCJ operation replaces an edge pair $\{\{a, b\}, \{c, d\}\}$ of the graph by $\{\{a, c\}, \{b, d\}\}$ or $\{\{a, d\}, \{b, c\}\}$. This edge edit operation on a graph is called a *2-break*.

This paper establishes a general framework for weighting rearrangements. The results are based on the problem of transforming one labeled graph into another through a scenario of operations, each weighted by an arbitrary function φ . The problem, called φ -MINIMUM COST PARSIMONIOUS SCENARIO (or φ -MCPS), asks for a scenario with a minimum number of 2-breaks, such that the sum of the costs for the operations is minimized.

Applications of our framework

While our framework is general, we use it to render two previous studies more practical. The first study is our work relating the likelihood of rearrangement breakpoints to the physical proximity in the nucleus [11]. This work is based on the hypothesis that two breakpoints could be confused when they are physically close. The model in this study labels the breakpoint graph edges (corresponding to intergenic regions) with fixed “colors”, and the cost of a DCJ has a weight of one if the labels are different and a weight of zero if they are the same. Using that cost function, we colored intergenic regions by grouping them according to their physical proximity, as inferred by Hi-C data. Although this technique of grouping proved to make biological sense [10, 12], it is far from ideal since much of the information given by the Hi-C data is lost in the labeling, and it is not immediately clear how to best compute the grouping. Our results here bypass the complexity of grouping by allowing each DCJ to be weighted by the values taken directly from the Hi-C contact maps. We give an algorithm for φ -MCPS on a breakpoint graph with an arbitrary φ and fixed edge labels, that runs in $O(n^5)$ time in the worst case but has better parameterized complexity in practice (see Example 1). We give in “Practical matters” section other reasons why the running times for this algorithm should remain practical.

The second study that we improve is that of Bulteau et al. [7]. Their biological constraint is based on the number of nucleotides in the intergenic regions containing breakpoints; they compute parsimonious scenarios that minimize the number of nucleotides inserted and deleted in intergenic regions. Their algorithm is restricted to instances where the breakpoint graph has only cycles (and no paths—sometimes referred to as *co-tailed* genomes). Using their $O(n \log n)$ algorithm, our

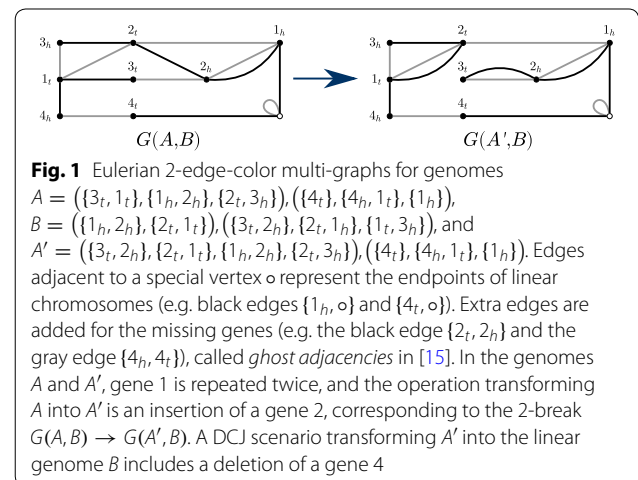
framework gives an $O(n^3)$ algorithm on any breakpoint graph (see Example 3).

This is an example of how our framework simplifies algorithm design on weighted DCJs. For a weight function adhering to our general criteria of “Cost-constrained 2-breaks” section, future algorithm designers now need only to concentrate on developing an efficient algorithm that works on a single cycle of a breakpoint graph. Thanks to Theorem 3, they will get a polynomial time algorithm that works on a general instance for free. “ α -approximation for ϕ -MCPS” section shows that the same is true for approximation algorithms.

This paper is based on general results we obtain on weighted transformations of edge-labeled multi-graphs. The permitted transformations can change the connectivity of the graph through a 2-break, or change the edge labels, or both. This model not only proves to be powerful enough to subsume the previously mentioned results, but also offers other advantages. It is flexible enough so that DCJ costs can be based on the labels of edges in the breakpoint graph, or on the labels of the vertices, or a combination of both. Also, since single-gene insertions and deletions can be represented as “ghost” adjacencies [15], all of this paper applies to genomes where genes could be missing in one genome or the other. Most results can be applied to genomes with duplicate genes (as depicted in Fig. 1).

Our model and general results

The foundation of this paper is a model for cost-constraining scenarios of degree preserving graph transformations, called 2-breaks, that are also known as edge swaps, switches, rewirings, or flips [16]. A 2-break transforms a graph by replacing two edges $\{u, v\}$ and $\{q, s\}$ by $\{u, q\}$ and $\{v, s\}$. These transformations have been studied not only in a restricted setting for genome



rearrangement [14, 17] and sorting strings by mathematical transpositions [18, 19], but also in the more general settings of generating random networks [16] and network design [20, 21].

Our results are about the transformation of an arbitrary multi-graph G into another one H having the same degree sequence. We find it convenient to reason in a setting, where we are given an Eulerian 2-edge-colored multi-graph with black and gray edges, the black edges being from G and the gray from H . We transform the connectivity of the black edges into the connectivity of the gray edges using a sequence of 2-breaks. Therefore, whenever we use the word *graph*, *path* (respectively *cycle*), we are referring to an Eulerian 2-edge-colored multi-graph, a path (respectively cycle) that alternates between black and gray edges. Naturally, a *cycle decomposition* of a graph is a partition of the edges of an Eulerian 2-edge-colored multi-graph into a set of alternating cycles. A *breakpoint graph* is a graph with a vertex for each gene extremity—each incident to exactly one gray and one black edge—along with one chromosome endpoint vertex \circ that could have degree as high as $2n$ (see Fig. 2). “DCJ scenarios for genomes and breakpoint graphs” section introduces the breakpoint graph in detail, and defines the Double Cut and Join (DCJ) model.

Our model for weighting 2-breaks is primarily based on a graph labeling, a set \mathcal{O} of valid operations, and a weight function $\varphi : \mathcal{O} \rightarrow \mathbb{R}_+$. Roughly speaking, a labeled input graph can be transformed through a series of operations in \mathcal{O} , where an operation can change the connectivity of the black edges of the graph, and/or change the labels of the edges. Any weight function φ defines an optimization problem φ -MCPS, which asks for a scenario that minimizes the total weight of the operations. This model subsumes many previously studied weighted DCJ models, as described in “Examples of the cost-constrained DCJ problems in the literature” section.

The spine of our results is built from successive theorems that speak to the decomposability into subproblems of a φ -MCPS instance. Lemma 3 shows that a parsimonious scenario of 2-breaks transforming the black edges into the gray implies a MAXIMUM ALTERNATING EDGE-DISJOINT CYCLE DECOMPOSITION (or MAECD) [22]. Theorem 1 says that an optimal solution to φ -MCPS can

be found using solutions to the MAECD problem, so that if φ -MCPS can be solved on a simple alternating cycle, then it can be solved on any instance. Theorem 2 says that an optimal solution to φ -MCPS on a simple alternating cycle can be found using a solution to the φ -MCPS problem on what we call a *circle*, that is, an alternating cycle that does not visit the same vertex twice (see Fig. 4).

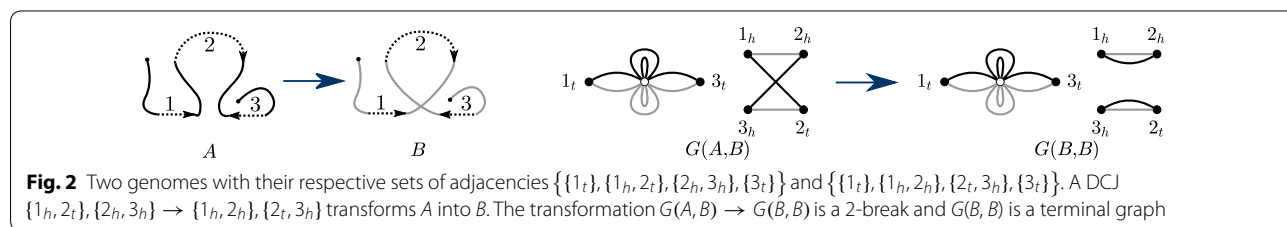
Under the common genome model, where each gene occurs exactly once in each genome, a relationship exists between parsimonious DCJ scenarios and solutions to MAECD on a breakpoint graph [14, 23]. We exploit this link in “ φ -MCPS for a breakpoint graph” section. Theorem 3 ties everything together; an amortized analysis shows that, given an $O(r^t)$ algorithm for computing φ -MCPS on a circle with r edges, φ -MCPS can be calculated on a breakpoint graph in $O(n^{t+1})$ time.

Under a more general genome model, that allows for changes in copy numbers of genes (e.g. insertions, deletions, and duplications), the spine of our results still holds due to the convenient representation of missing genes as *ghost adjacencies* in an Eulerian 2-edge-colored multi-graph [15] (see Fig. 1). All of our results hold for pairs of genomes with non-duplicated genes, but unequal gene content. Indeed, a breakpoint graph (i.e. graph with limited degree for most nodes) can still represent the pair of genomes in this case.

Caprara proved that MAECD is NP-Hard for Eulerian 2-edge-colored multi-graphs where each vertex is incident to at most two gray and two black edges (which is the case when there are two copies of each gene) [22]. We present a simple integer linear program (or ILP) that solves φ -MCPS for these types of graphs, given a method to solve φ -MCPS on a circle. This ILP is likely to be unwieldy in general, since the number of variables is exponential in the number of simple alternating cycles. In the case of breakpoint graphs on specific genomes, this may not always be intractable, as the number of duplicate genes may be limited. See “Practical matters” section for a discussion of these practical matters.

DCJ scenarios for genomes and breakpoint graphs

A *genome* consists of *chromosomes* that are linear or circular orders of genes separated by potential *breakpoint* regions. In Fig. 2 the tail of an arrow represents the *tail*



extremity, and the head of an arrow represents the *head extremity* of a gene. We can represent a genome by a set of *adjacencies* between the gene extremities. An *adjacency* is either *internal*: an unordered pair of the extremities that are adjacent on a chromosome, or *external*: a single extremity adjacent to one of the two ends of a linear chromosome. In what follows we will suppose that two genomes A and B are partitioned into n genes each occurring exactly once in each genome, and our goal will be to transform A into B using a sequence of DCJs.

Definition 1 (*double cut and join*) A DCJ cuts one or two breakpoint regions and joins the resulting ends of the chromosomes back in one of the four following ways: $\{a, b\}, \{c, d\} \rightarrow \{a, c\}, \{b, d\}$; $\{a, b\}, \{c\} \rightarrow \{a, c\}, \{b\}$; $\{a, b\} \rightarrow \{a\}, \{b\}$; and $\{a\}, \{b\} \rightarrow \{a, b\}$.

We represent the pairs of the genomes with a help of a breakpoint graph [13, 17].

Definition 2 (*breakpoint graph*) $G(A, B)$ for genomes A and B is a 2-edge-colored Eulerian undirected multi-graph. V consists of $2n$ gene extremities and an additional vertex \circ . For every internal adjacency $\{a, b\} \in A$ (respectively $\{a, b\} \in B$) there is a black (respectively gray) edge $\{a, b\}$ in $G(A, B)$ and for every external adjacency $\{a\} \in A$ (respectively $\{a\} \in B$) there is a black (respectively gray) edge $\{a, \circ\}$ in $G(A, B)$. There is a number of black and gray loops $\{\circ, \circ\}$ ensuring that $d^b(G(A, B), \circ) = d^g(G(A, B), \circ) = 2n$.

2-Break scenarios for 2-edge-colored graphs

In this paper a *graph* is an Eulerian 2-edge-colored undirected multi-graph with edges colored black or gray as in Fig. 1. A graph with equal multi-sets of black and gray edges is called *terminal*, and our goal is to transform a given graph into a terminal one using 2-breaks.

Definition 3 (*2-break scenario*) A 2-break replaces two black edges $\{x_1, x_2\}$ and $\{x_3, x_4\}$ by either $\{x_1, x_3\}$ and $\{x_2, x_4\}$ or $\{x_1, x_4\}$ and $\{x_2, x_3\}$. A 2-break *scenario* of length m is a sequence of m 2-breaks transforming a graph into a terminal one.

Definition 4 (*Eulerian graph and alternating cycle*) G is *Eulerian* if every vertex has equal black and gray degrees. A cycle is *alternating* if it is Eulerian. All use of the word *cycle* in this paper will be synonymous with alternating cycle.

Define a MAXIMUM ALTERNATING EDGE-DISJOINT CYCLE DECOMPOSITION (MAECD) of a graph G as a

decomposition of G into a maximum number of edge-disjoint alternating cycles. Denote the size of a MAECD of G by $c(G)$ and the number of its black edges by $e(G)$. We make a distinction between simple cycles and circles (see Fig. 4 to see a simple cycle that is not a circle).

Definition 5 (*simple cycle and circle*) A graph G is a *simple cycle* if the size of a MAECD, $c(G) = 1$. If in addition to that the black and gray degrees $deg^b(G, v)$ and $deg^g(G, v)$ are equal to 1 for every vertex v , then G is called a *circle*.

Parsimonious 2-break scenarios

The problem of finding a minimum length (or *parsimonious*) 2-break scenario was treated in several unrelated settings using different terminology. Lemma 1 proven in “Proofs” section was treated in [20] where the authors also showed that finding a minimum length 2-break scenario is NP-hard due to the NP-hardness of finding a MAECD of a graph and provided a $7/4$ -approximation algorithm for finding this length. A variant of the problem for Eulerian digraphs where all the gray edges are loops was solved in [24].

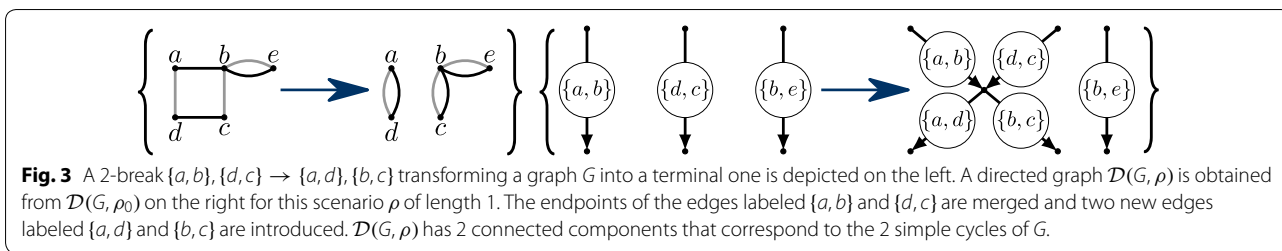
Lemma 1 (*Bienstock and Günlük in [20]*) *The minimum length of a 2-break scenario on a graph G is $d_{2b}(G) = e(G) - c(G)$.*

Since finding a MAECD for a breakpoint graph is easy, Lemma 1 leads to a linear time algorithm for finding a parsimonious DCJ scenario [13]. The algorithm is based on Lemma 2 proven in “Proofs” section.

Lemma 2 (*Yancopoulos et al. in [13]*) *The minimum length of a DCJ scenario transforming genome A into B is equal to $d_{2b}(G(A, B)) = e(G(A, B)) - c(G(A, B))$.*

Decomposition of a 2-break scenario

In this section we will show how a 2-break scenario ρ of length m can be partitioned into subscenarios ρ^1, \dots, ρ^k and G can be decomposed into edge-disjoint Eulerian subgraphs H^1, \dots, H^k where ρ^i is a scenario for H^i , and $k \geq e(G) - m$. We will use this decomposition in “ φ -MCPS for a graph” section to show that φ -MCPS on a graph can be solved by solving φ -MCPS on its simple cycles. For a graph G and a 2-break scenario ρ we define a directed 1-edge-colored edge-labeled graph $\mathcal{D}(G, \rho)$, akin to the *trajectory graph* introduced by Shao et al. [25]. Denote the sequence of the first l 2-breaks of ρ by ρ_l and the graph obtained from G after these 2-breaks by G_l . Define $\mathcal{D}(G, \rho_0)$ in the following way: for each black edge e of G we have two new vertices connected by a directed



edge labeled by e (see Fig. 3). For the l -th 2-break of ρ , $\{x_1, x_2\}, \{x_3, x_4\} \rightarrow \{x_1, x_3\}, \{x_2, x_4\}$, merge the endpoints of the edges labeled $\{x_1, x_2\}$ and $\{x_3, x_4\}$ in $\mathcal{D}(G, \rho_{l-1})$. Proceed by adding two new vertices to $\mathcal{D}(G, \rho_{l-1})$ and two edges labeled $\{x_1, x_3\}$ and $\{x_2, x_4\}$ from the merged vertex to the newly added ones to obtain $\mathcal{D}(G, \rho_l)$. Continue until $\mathcal{D}(G, \rho_m)$ is obtained, where m is the length of ρ , and denote it by $\mathcal{D}(G, \rho)$.

Shao et al. [25] characterize the connected components of a trajectory graph for a parsimonious scenario. Using similar techniques we prove the following lemma in “Proofs” section.

Lemma 3 *If $\mathcal{D}(G, \rho)$ has k connected components then ρ can be partitioned into k subscenarios ρ^i and G can be partitioned into k edge-disjoint Eulerian subgraphs H^i in such a way that ρ^i is a scenario for H^i for every $i \in \{1, \dots, k\}$. If ρ is parsimonious, then $k = c(G)$ and $C(\rho) = \{H^1, \dots, H^k\}$ is a MAECD of G .*

Cost-constrained 2-breaks

In this section we outline our model for assigning costs to 2-breaks. We associate labels to both vertices and edges of a graph, and then describe a set \mathcal{O} of *valid operations* of 2-breaks on labeled edges and edge-label changes. Our cost function is defined on \mathcal{O} . This model generalizes the labeled DCJ problems of [7, 11].

We will use letters u, v, q, s to denote vertices, letters a, b, c, d to denote vertex labels and x, y, z, t to denote edge labels. Given an alphabet of vertex labels Σ_V and one of edge labels Σ_E , fix a subset \mathcal{O} containing a set of tuples

- $((\{a, b\}, x); (\{a, b\}, y))$ (called edge-label changes) and
- $((\{a, b\}, x), (\{c, d\}, y); (\{a, c\}, z), (\{b, d\}, t))$ (called 2-breaks on labels)

for $a, b, c, d \in \Sigma_V$ and $x, y, z, t \in \Sigma_E$.

Take a graph $G = (V, E)$, and its labeling $\lambda = (\lambda_V, \lambda_E)$ with $\lambda_V : V \rightarrow \Sigma_V$ and $\lambda_E : E \rightarrow \Sigma_E$. If \mathcal{O} contains an edge-label change $((\{a, b\}, x); (\{a, b\}, y))$ and (G, λ) contains an edge $\{u, v\}$ labeled x with vertices u and v labeled

a and b , then the label of this edge can be changed into y . We call such a transformation of (G, λ) an \mathcal{O} -change and denote it $((\{u, v\}, x) \rightarrow (\{u, v\}, y))$.

If \mathcal{O} contains a 2-break on labels $((\{a, b\}, x), (\{c, d\}, y); (\{a, c\}, z), (\{b, d\}, t))$ and (G, λ) contains two edges $\{u, v\}$ and $\{q, s\}$ labeled x and y respectively with vertices u, v and q, s labeled a, b and c, d , then a 2-break $\{u, v\}, \{q, s\} \rightarrow \{u, q\}, \{v, s\}$ can be performed on G with the labels of the new edges being z and t . We call such a transformation of (G, λ) an \mathcal{O} -break and denote it $((\{u, v\}, x), (\{q, s\}, y) \rightarrow (\{u, q\}, z), (\{v, s\}, t))$.

An \mathcal{O} -scenario $\rho_{\mathcal{O}}$ for (G, λ) is a sequence of \mathcal{O} -changes and \mathcal{O} -breaks transforming (G, λ) into $(\bar{G}, \bar{\lambda})$ such that \bar{G} is terminal and its multi-sets of black and gray labeled edges are equal. The number of \mathcal{O} -breaks in $\rho_{\mathcal{O}}$ will be called the *2-break-length* of the scenario. If a $\rho_{\mathcal{O}}$ exists for (G, λ) , then $d_{\mathcal{O}b}(G, \lambda)$ denotes the minimum 2-break-length of an \mathcal{O} -scenario.

An \mathcal{O} -scenario does not necessarily exist for a given (G, λ) , however if it exists, then the inequality $d_{\mathcal{O}b}(G, \lambda) \geq d_{2b}(G)$ holds, where $d_{2b}(G)$ is the minimum length of a 2-break scenario on a graph G . In this paper we deal with the sets \mathcal{O} that have the necessary operations to parsimoniously transform (G, λ) into $(\bar{G}, \bar{\lambda})$. We call these sets *p-sufficient*.

Definition 6 (*p-sufficient \mathcal{O} for (G, λ)*) A set \mathcal{O} is *parsimonious-sufficient* or *p-sufficient* for (G, λ) if we have $d_{\mathcal{O}b}(G, \lambda) = d_{2b}(G)$.

The cost function that we consider is $\varphi : \mathcal{O} \rightarrow \mathbb{R}_+$. The cost of an \mathcal{O} -scenario is the sum of the costs of its constituent operations. If \mathcal{O} is p-sufficient for (G, λ) , then $MCPS_{\varphi}(G, \lambda)$ is the minimum cost of an \mathcal{O} -scenario of the 2-break-length equal to $d_{2b}(G)$, otherwise $MCPS_{\varphi}(G, \lambda)$ is ∞ . We consider the following problem:

Problem 1 (φ -MINIMUM COST PARSIMONIOUS SCENARIO or φ -MCPS)

INPUT: A graph G , and its labeling λ .
 OUTPUT: $MCPS_{\varphi}(G, \lambda)$.

Examples of the cost-constrained DCJ problems in the literature

Example 1 (MINIMUM LOCAL PARSIMONIOUS SCENARIO) In [11] we supposed the adjacencies of genome A to be partitioned into spatial regions represented by different colors. We then developed a polynomial time algorithm for finding a parsimonious DCJ scenario minimizing the number of rearrangements whose breakpoints appear in different regions. The problem as was stated in [11] differs slightly from φ -MCPS, since in that study we do not have colors for the adjacencies of genome B . We can bridge this gap as follows.

Edge labels $\Sigma_E = \Sigma_c \cup \{\tau\}$ are the colors representing the different spatial regions of a genome plus an additional terminal label τ . There is a single vertex label $\Sigma_V = \{a\}$. \mathcal{O} contains 2-breaks on labels $((\{a, a\}, x), (\{a, a\}, y); (\{a, a\}, x), (\{a, a\}, y))$ for $x, y \in \Sigma_c$, and edge-label changes $((\{a, a\}, x); (\{a, a\}, \tau))$ for $x \in \Sigma_c$. The cost φ_c of a 2-break on labels in \mathcal{O} is 0 if the 2 labels being replaced are equal and 1 otherwise. The cost of a edge-label change is 0.

In [11] we presented an $O(n^4)$ time algorithm solving φ_c -MCPS for a labeled breakpoint graph with the gray edges labeled by τ . In [12] we demonstrated that finding a minimum cost \mathcal{O} scenario for such a breakpoint graph, when the parsimonious criteria is disregarded, is NP-hard. We proposed an algorithm that is exponential in the number of colors but not in the number of genes.

In “ φ_f -MCPS for a circle with fixed labels” section we use the same \mathcal{O} , fix a symmetric function $\Phi : \Sigma^2 \rightarrow \mathbb{R}_+$, and define $\varphi_f((\{a, a\}, x), (\{a, a\}, y); (\{a, a\}, x), (\{a, a\}, y)) = \Phi(x, y)$. This drastically enhances the model introduced in [11] as now rearrangements whose breakpoints appear in the same region can have non-zero costs. In “ φ_f -MCPS for a breakpoint graph” section we provide an $O(n^5)$ time algorithm solving the generalized problem of φ_f -MCPS for a labeled breakpoint graph.

Example 2 (DCJ WEIGHTED BY HI-C) In [10] we weighted each DCJ by the value taken directly from the Hi-C contact map. In this model every intergenic region of genome A gets assigned an interval corresponding to its genomic coordinates on a chromosome. The *weight* of a DCJ acting on two intergenic regions is then equal to the average Hi-C value for their corresponding intervals. In [10] we presented an algorithm greedily maximizing the weight of a parsimonious scenario and found that the

obtained weight is significantly higher than the weight of a random parsimonious scenario.

Edge labels are the genomic intervals corresponding to the intergenic regions of a genome A plus an additional terminal label. There is a single vertex label $\Sigma_V = \{a\}$. \mathcal{O} stays as in Example 1. $\Phi_{HiC}(x, y)$ on two genomic intervals is their average Hi-C value. The problem that maximizes Hi-C values can be easily transformed into a minimization problem by setting the cost of a 2-break on labels $((\{a, a\}, x), (\{a, a\}, y); (\{a, a\}, x), (\{a, a\}, y))$ to $\Phi_{\max} - \Phi_{HiC}(x, y)$, where Φ_{\max} is the maximum $\Phi_{HiC}(x, y)$ over all x, y .

In [10] the optimality of the proposed greedy algorithm was not discussed, but our work presented in “ φ_f -MCPS for a circle with fixed labels” section of this paper provides us with a polynomial time algorithm for solving this problem exactly.

Example 3 (SORTING BY wDCJs AND INDELS IN INTERGENES) Bulteau et al. [7] introduced a problem where adjacencies of genomes are labeled with their genetic length (number of nucleotides). A *wDCJ* is a DCJ that preserves the sum of the genetic lengths of the adjacencies and an *indel* δ increases or decreases the genetic length of an adjacency by δ . The cost of a wDCJ is 0 and the cost of an indel δ is $|\delta|$. A scenario of wDCJs and indels for (G, λ) is said to be *valid* if its wDCJ-length is $d_{2b}(G)$. The paper presents an $O(n \log n)$ algorithm for finding a minimum cost scenario among the *valid* ones, for the genomes with circular chromosomes and n genes.

Translating this into our formalism yields the following φ -MCPS problem. Edge labels are the natural numbers, there is a single vertex label, and \mathcal{O} contains 2-breaks on labels $((\{a, a\}, w_1), (\{a, a\}, w_2); (\{a, a\}, w_3), (\{a, a\}, w_4))$ for $w_i \in \Sigma_E$ satisfying $w_1 + w_2 = w_3 + w_4$. \mathcal{O} also contains edge-label changes $((\{a, a\}, w_1); (\{a, a\}, w_2))$ for $w_i \in \Sigma$. \mathcal{O} is p-sufficient for any (G, λ) since G can be first transformed into a terminal graph using any parsimonious 2-break scenario and then its labels can be adjusted. The cost φ_l of a 2-break on labels is 0 and the cost φ_l of a edge-label change $((\{a, a\}, w_1); (\{a, a\}, w_2))$ is $|w_1 - w_2|$.

In [7] the authors presented an $O(r \log r)$ time algorithm for solving φ_l -MCPS on a circle with r vertices. Combining this algorithm with our results from “ φ_f -MCPS for a breakpoint graph” section gives an algorithm solving φ_l -MCPS in $O(n^3)$ time for a labeled breakpoint graph. The ILP defined in “ φ_f -MCPS for a graph” section solves φ_l -MCPS for any labeled graph.

Example 4 (wDCJ-DIST) Fertin et al. [6] treated a problem wDCJ-DIST where wDCJs without indels are allowed, and the sums of the genetic lengths of the adjacencies of two genomes are equal.

In this case we keep the same Σ_E, Σ_V and \mathcal{O} as in Example 3 except that the edge-label changes are excluded from \mathcal{O} . A labeled graph is said to be *balanced* if the sums of the labels of black and gray edges are equal. wDCJ-DIST is the problem of finding $d_{\mathcal{O}b}$ for a balanced graph whose connected components are circles. The authors show that wDCJ-DIST is strongly NP-complete. However they also prove that $d_{\mathcal{O}b}(O, \lambda) = d_{2b}(O)$ for a balanced circle O and that \mathcal{O} is p-sufficient for a graph whose connected components are balanced circles.

Example 5 Although ignored in the previous examples, the weighting of operations based on only the vertices is also possible under our framework. For example, take $\Sigma_E = \{\tau\}, \Sigma_V = \mathbb{N}, \mathcal{O}$ containing 2-breaks on labels $((\{a, b\}, \tau), (\{c, d\}, \tau); (\{a, c\}, \tau), (\{b, d\}, \tau))$ and any cost function $\varphi_v : \mathcal{O} \rightarrow \mathbb{R}_+$. The costs of the 2-breaks on labels in \mathcal{O} could be a function of the genomic coordinates of the participating gene extremities.

Note that the set \mathcal{O} is implicit, rather than explicit. In Example 3, \mathcal{O} would be too large to represent explicitly since every pair of genetic lengths for every pair of edges would exist. For all practical uses that we know of to date, membership in \mathcal{O} can be computed in constant time.

Suppose that there exists a MAECD C of G consisting of the simple cycles for which \mathcal{O} is p-sufficient. For every $S \in C$ take an \mathcal{O} -scenario $\rho_{\mathcal{O}}^S$ of cost $MCPS_{\varphi}(S, \lambda^S)$ and 2-break-length $d_{2b}(S)$. By performing these scenarios one after another we obtain an \mathcal{O} -scenario $\rho_{\mathcal{O}}$ for (G, λ) of 2-break-length $\sum_{S \in C} d_{2b}(S) = d_{2b}(G)$ and of cost $\sum_{S \in C} MCPS_{\varphi}(S, \lambda^S)$. This means that $MCPS_{\varphi}(G, \lambda) \leq \sum_{S \in C} MCPS_{\varphi}(S, \lambda^S)$.

On the other hand, suppose that \mathcal{O} is p-sufficient for (G, λ) and take an \mathcal{O} -scenario $\rho_{\mathcal{O}}$ for (G, λ) of length $d_{2b}(G)$. For ρ , a 2-break scenario obtained from $\rho_{\mathcal{O}}$ when the labels are neglected, a decomposition $C(\rho)$ corresponding to ρ is a MAECD of G due to Lemma 3. A subsequence $\rho_{\mathcal{O}}^S$ of $\rho_{\mathcal{O}}$, consisting of the operations acting on the edges of a cycle $S \in C(\rho)$, is an \mathcal{O} -scenario for (S, λ^S) of 2-break-length $d_{2b}(S)$. A sequence of operations $\hat{\rho}_{\mathcal{O}}$ obtained by performing the subsequences $\rho_{\mathcal{O}}^S$ one after another for each $S \in C(\rho)$ is an \mathcal{O} -scenario for (G, λ) . By construction the 2-break-length of $\hat{\rho}_{\mathcal{O}}$ is equal to the 2-break-length of $\rho_{\mathcal{O}}$. The costs of $\rho_{\mathcal{O}}$ and $\hat{\rho}_{\mathcal{O}}$ are also equal, as they consist of exactly the same operations that are performed in different orders, thus the cost of $\rho_{\mathcal{O}}$ is greater or equal to $\sum_{S \in C(\rho)} MCPS_{\varphi}(S, \lambda^S) \geq \min\{\sum_{S \in C} MCPS_{\varphi}(S, \lambda^S) \mid C \text{ is a MAECD of } G\}$. \square

Take the set \mathcal{S} of simple labeled cycles of (G, λ) . If one can solve φ -MCPS for every $S \in \mathcal{S}$, then Theorem 1 provides a straightforward way to solve φ -MCPS for (G, λ) as a set packing problem. First compute $c(G)$ by solving the ILP in the left column. Then proceed by solving the other ILP to compute $MCPS_{\varphi}(G, \lambda)$.

$$\begin{aligned} & \text{Maximize } \sum_{S \in \mathcal{S}} x_S \\ & \text{Subject to } \sum_{S: e \in S} x_S \leq 1 \text{ for each edge } e \text{ of } G \\ & \text{and } x_S \in \{0, 1\} \text{ for simple cycle } S \in \mathcal{S}. \end{aligned}$$

$$\begin{aligned} & \text{Minimize } \sum_{S \in \mathcal{S}} x_S MCPS_{\varphi}(S, \lambda^S) \\ & \text{Subject to } \sum_{S: e \in S} x_S \leq 1 \text{ for each edge } e \text{ of } G, \\ & \sum_{S \in \mathcal{S}} x_S = c(G) \\ & \text{and } x_S \in \{0, 1\} \text{ for simple cycle } S \in \mathcal{S}. \end{aligned}$$

φ -MCPS for a graph

Theorem 1 Denote the φ -cost of a MAECD as the sum of the $MCPS_{\varphi}$ on its cycles. $MCPS_{\varphi}$ for a graph is equal to the minimum φ -cost of its MAECD.

Proof For a cycle S of a labeled graph (G, λ) , λ^S denotes the labeling of S according to λ . We suppose that $\min(\emptyset) = \infty$ and prove the following:

$$\begin{aligned} & MCPS_{\varphi}(G, \lambda) \\ & = \min \left\{ \sum_{S \in C} MCPS_{\varphi}(S, \lambda^S) \mid C \text{ is a MAECD of } G \right\}. \end{aligned}$$

There exists an algorithm efficiently listing all the simple cycles of an undirected 1-edge-colored graph [26], however we are unaware of a similar result for the 2-edge-colored graphs. Computing $c(G)$ is an APX-hard problem [27] and the size of \mathcal{S} may be exponential in the size of G , which might make these ILPs intractable in general. For graphs representing genomes with duplicate genes, the number of simple cycles can grow exponentially as a function of the number of duplicate genes. For breakpoint graphs, however, the number grows quadratically and $c(G)$ can be found in linear time.

φ -MCPS for a simple cycle

The decomposition theorem of “ φ -MCPS for a graph” section reduces the computation of φ -MCPS on a graph to the computation of φ -MCPS on a simple alternating cycle. In this section we further decompose the problem into simpler versions of cycles, called circles, which are alternating cycles that contain a vertex only once.

Denote $deg_2(G)$ for a graph G as the number of vertices with black and gray degree equal to two. It is easy to check that $deg^b(S, v) = deg^g(S, v) \leq 2$ for any vertex v of a simple cycle S . If $deg_2(S) = 0$, then S is a circle. See the first column of Fig. 4 for examples of simple cycles that are not circles.

Take a simple labeled cycle (S, λ) and denote S_0 as $\{(S, \lambda)\}$. Choose a vertex v of degree two in S and replace it by two vertices $v_1, v_2 \notin V$ labeled by the same label as v . If v is incident to a gray loop, then split it into two vertices v_1 and v_2 , as depicted on the top row of Fig. 4, to obtain a set S_1 consisting of a single simple cycle. Otherwise split it into two vertices, as depicted on the bottom row of Fig. 4, to obtain a set S_1 consisting of two simple cycles.

Simple labeled cycles in S_1 share the same set of vertices of degree two. Choose such a vertex and split it simultaneously in all the cycles in S_1 as previously to obtain a set S_2 of at most 4 simple labeled cycles sharing the same set of vertices and the same multi-set of labeled black edges. Continue this procedure until the set $circ(S, \lambda) = S_{deg_2(S)}$ of the labeled circles is obtained.

Theorem 2 $MCPS_\varphi$ for a simple cycle (S, λ) is equal to the minimum of the $MCPS_\varphi$ among the circles in $circ(S, \lambda)$.

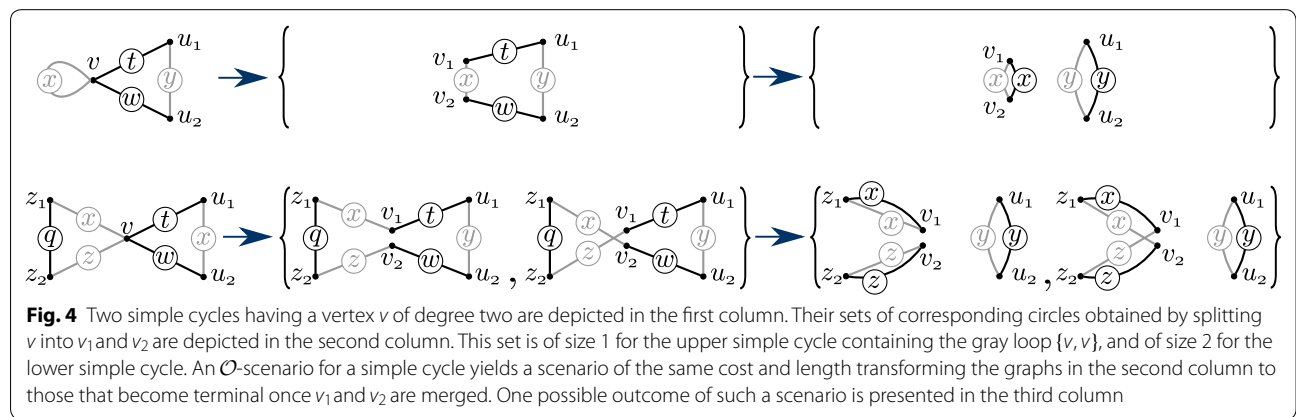
Proof First we prove that $MCPS_\varphi(S, \lambda) = \min\{MCPS_\varphi(H, \mu) \mid (H, \mu) \in S_1\}$. Labeled graphs in S_1 are obtained by splitting a vertex v of degree 2 into vertices v_1 and v_2 . For a labeled graph (H, μ) on vertices $V \cup \{v_1, v_2\} \setminus \{v\}$ denote $r_g(H, \mu)$ as the labeled graph obtained from (H, μ) by reversing the split, that is, by merging the vertices v_1 and v_2 into v .

Choose $(\hat{S}, \hat{\lambda}) \in S_1$. By construction $r_g(\hat{S}, \hat{\lambda}) = (S, \lambda)$. Denote $r_v(v_1) = r_v(v_2) = v$, and $r_v(u) = u$ for $u \in V$. For an edge f of $(\hat{S}, \hat{\lambda})$ joining vertices q and s , the edge $r_e(f) = \{r_v(q), r_v(s)\}$ is present in (S, λ) and has the same label as f . r_e defines a bijection between the labeled edges of (S, λ) and $(\hat{S}, \hat{\lambda})$ and thus between \mathcal{O} operations on these graphs. This means that an operation in \mathcal{O} transforming $(\hat{S}, \hat{\lambda})$ into some $(\hat{S}', \hat{\lambda}')$ transforms (S, λ) into $r_g(\hat{S}', \hat{\lambda}')$, and an operation in \mathcal{O} transforming (S, λ) into some (S', λ') transforms $(\hat{S}, \hat{\lambda})$ into $(\hat{S}', \hat{\lambda}')$ such that $r_g(\hat{S}', \hat{\lambda}') = (S', \lambda')$.

Thus for an \mathcal{O} -scenario of $(\hat{S}, \hat{\lambda})$ there exists an \mathcal{O} -scenario of the same φ cost and the same 2-break-length for (S, λ) . On the other hand, an \mathcal{O} -scenario for (S, λ) provides us with a sequence ρ of \mathcal{O} operations of the same φ cost and the same 2-break-length transforming $(\hat{S}, \hat{\lambda})$ into $(\bar{S}, \bar{\lambda})$ for which $r_g(\bar{S}, \bar{\lambda})$ is a terminal graph.

If S_1 is of size 1, then there is a single choice for $(\bar{S}, \bar{\lambda})$ (see the right upper corner of Fig. 4) and it is itself terminal. If S_1 is of size 2, then there are two options for $(\bar{S}, \bar{\lambda})$ (see the right bottom corner of Fig. 4). Either $(\bar{S}, \bar{\lambda})$ is already terminal, or the sequence ρ of \mathcal{O} operations transforming $(\hat{S}, \hat{\lambda})$ into $(\bar{S}, \bar{\lambda})$ transforms the second graph in S_1 into a terminal one.

Now we prove that $MCPS_\varphi(S, \lambda) = \min\{MCPS_\varphi(O, \lambda) \mid (O, \lambda) \in circ(S, \lambda)\}$, which is clearly true for $deg_2(S) = 0$. Suppose this to be true for $deg_2(S) < t$. We prove it for $deg_2(S) = t$ by induction. For $(\hat{S}, \hat{\lambda}) \in S_1$ one has $deg_2(\hat{S}) = t - 1$, so using the inductive hypothesis we have that $MCPS_\varphi(\hat{S}, \hat{\lambda})$ is equal to $\min\{MCPS_\varphi(O, \lambda) \mid (O, \lambda) \in circ(\hat{S}, \hat{\lambda})\}$. We have already proven that $MCPS_\varphi(S, \lambda) = \min\{MCPS_\varphi(H, \mu) \mid (H, \mu) \in S_1\}$, and by construction we know that $circ(S, \lambda) = \cup_{(H, \mu) \in S_1} circ(H, \mu)$. These combine to imply that the theorem is true for $deg_2(S) = t$. \square



φ -MCPS for a breakpoint graph

In this section we suppose that there exists an algorithm for computing $MCPS_\varphi$ on a labeled circle (e.g. the algorithm of “ φ_f -MCPS for a circle with fixed labels” section). Using this algorithm as a subroutine we will construct an algorithm for finding $MCPS_\varphi$ for a labeled breakpoint graph. This is a generalization of the work first presented in [11].

Take genomes A and B partitioned into n genes where each gene occurs exactly once in each genome, and a labeling λ of a breakpoint graph $G(A, B)$. For all the vertices $v \neq \circ$ we have $deg^g(G(A, B), v) = deg^b(G(A, B), v) = 1$. Thus, if there is a circle in $G(A, B)$ containing an edge then this circle is the only simple cycle containing this edge. This means that every MAECD of $G(A, B)$ includes all of its circles. These set aside we are left with $G(A, B)'$, which is a union of alternating paths starting and ending at \circ with end edges of the same color. If this color is black we call the path AA , and BB otherwise.

We proceed by constructing a complete weighted bipartite graph H having the AA and BB paths of $G(A, B)'$ as vertices. Every simple cycle of $G(A, B)'$ is a union of an AA path and a BB path. To each edge joining these paths in H we assign weight equal to $MCPS_\varphi$ for a union of these paths. A MAECD of $G(A, B)'$ corresponds to a maximum matching of H and every such matching corresponds to a MAECD of $G(A, B)'$. Denote λ' as the labeling of $G(A, B)'$ according to λ . Using Theorem 1 we obtain that $MCPS_\varphi(G(A, B)', \lambda')$ is equal to the minimum weight of a maximum matching of H . There is an equal number p of AA and BB paths. Let P denote the total number of edges in $G(A, B)'$. Using this notation we obtain the following lemma proven in “Proofs” section.

Lemma 4 *For a function f and an $O(f(r))$ time algorithm for φ -MCPS on a labeled circle on r vertices, there exists an $O(p^2f(P) + p^3 + f(n))$ time algorithm for φ -MCPS on a labeled breakpoint graph. If $f(r) = O(r^t)$ for some constant $t \geq 1$, then φ -MCPS on a labeled breakpoint graph can be solved in $O(pP^t + p^3 + n^t)$ time.*

Both p and P are $O(n)$, thus Lemma 4 leads to the following theorem.

Theorem 3 *Given a constant $t \geq 2$ and an $O(r^t)$ time algorithm for φ -MCPS on a labeled circle on r vertices, φ -MCPS on a labeled breakpoint graph can be solved in $O(n^{t+1})$ time.*

Corollary 1 *Using the $O(r^4)$ algorithm from “ φ_f -MCPS for a circle with fixed labels” section we obtain an $O(n^5)$ algorithm for solving φ_f -MCPS on a labeled breakpoint graph with fixed labels.*

Corollary 2 *Using the $O(r \log r)$ algorithm from [7] for the SORTING BY WDCs AND INDELS IN INTERGENES problem on a circle (see Example 3), we obtain an $O(n^3)$ algorithm for solving the problem on a breakpoint graph.*

α -approximation for φ -MCPS

Theorems 1 and 2 demonstrate how φ -MCPS for any labeled graph can be solved if one is able to solve φ -MCPS for a labeled circle. This is exploited in Theorem 3 to solve φ -MCPS for a breakpoint graph. Analogous results proven in “Proofs” section hold if instead of an exact algorithm one has an α -approximation for φ -MCPS for a labeled circle.

Lemma 5 *For a constant $t \geq 2$ and an $O(r^t)$ time α -approximation algorithm for φ -MCPS on a labeled circle on r vertices, there exists an $O(n^{t+1})$ time α -approximation algorithm for φ -MCPS on a labeled breakpoint graph.*

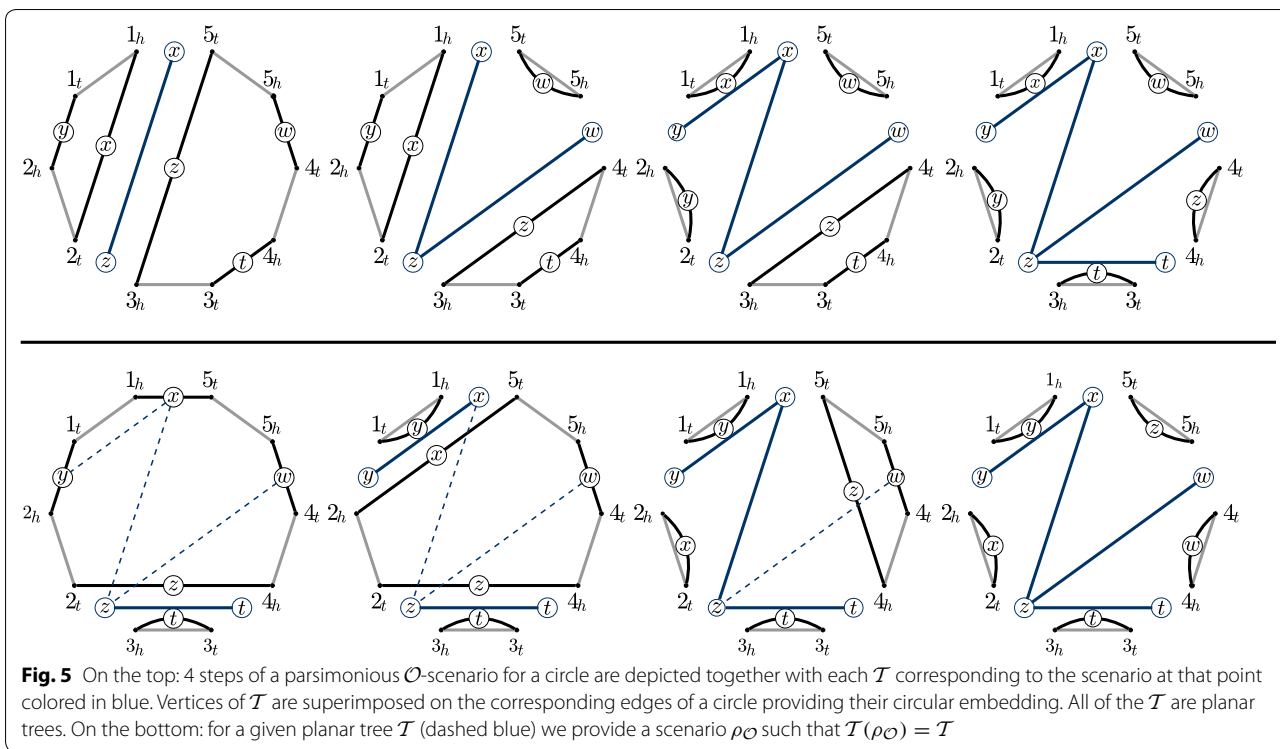
φ_f -MCPS for a circle with fixed labels

Here we define φ_f -MCPS, a particular instance of a φ -MCPS problem, and solve it for a circle. φ_f -MCPS generalizes our previous work presented in Examples 1 and 2. For a set $\Sigma_V = \{a\}$ of vertex labels and a set $\Sigma_E = \Sigma \cup \{\tau\}$ of edge labels, define a set \mathcal{O} consisting of 2-breaks on labels $((\{a, a\}, x), (\{a, a\}, y); (\{a, a\}, x), (\{a, a\}, y))$ for $x, y \in \Sigma$, and edge-label changes $((\{a, a\}, x); (\{a, a\}, \tau))$ for $x \in \Sigma$. Fix a symmetric function $\Phi : \Sigma^2 \rightarrow \mathbb{R}_+$ and define a φ_f cost of a 2-break on labels $((\{a, a\}, x), (\{a, a\}, y); (\{a, a\}, x), (\{a, a\}, y))$ to be $\Phi(x, y)$ and a φ_f cost of an edge-label change $((\{a, a\}, x); (\{a, a\}, \tau))$ to be 0. We will provide a polynomial time algorithm for φ_f -MCPS on a labeled circle with the gray edges labeled by a terminal label τ .

Without loss of generality we can suppose that all of the black edges of a circle have different labels; if two edges are labeled with the same label x , then we simply replace one of these labels with a new label \hat{x} and set $\hat{\Phi}(\hat{x}, y) = \Phi(x, y)$ and $\hat{\Phi}(y, z) = \Phi(y, z)$ for $y, z \in \Sigma$.

For a labeled circle having r black edges, define a set V_Σ of r vertices corresponding to their labels. For an \mathcal{O} -scenario $\rho_{\mathcal{O}}$ we define a 1-edge-colored undirected graph $\mathcal{T}(\rho_{\mathcal{O}})$ with vertices V_Σ and an edge $\{x, y\}$ for every \mathcal{O} -break in $\rho_{\mathcal{O}}$ replacing the black edges labeled with x and y (see Fig. 5). The cost of an edge $\{x, y\}$ is defined to be $\Phi(x, y)$ and the cost of a graph $\mathcal{T}(\rho_{\mathcal{O}})$ is the sum of the costs of its edges. The costs of $\rho_{\mathcal{O}}$ and $\mathcal{T}(\rho_{\mathcal{O}})$ are equal by construction.

Fix a circular embedding of V_Σ respecting the order of the black edges on the labeled circle (see Fig. 5). A graph with vertices V_Σ is said to be *planar on the circle* if none of its edges cross in this embedding. We prove Lemma 6 in “Proofs” section linking planar trees and parsimonious scenarios.



Lemma 6 *If $\rho_{\mathcal{O}}$ is a minimum 2-break-length \mathcal{O} -scenario for a labeled circle (\mathcal{O}, λ) , then $\mathcal{T}(\rho_{\mathcal{O}})$ is a planar tree on (\mathcal{O}, λ) . In addition to that, for a planar tree \mathcal{T} on (\mathcal{O}, λ) there exists an \mathcal{O} -scenario $\rho_{\mathcal{O}}$ such that $\mathcal{T}(\rho_{\mathcal{O}}) = \mathcal{T}$.*

Farnoud and Milenkovic in [19] provide a dynamic programming algorithm for finding a minimum cost planar tree on a circle. In “Proofs” section their proof for a following lemma is given which, together with Lemma 6, leads to Theorem 4.

Lemma 7 (Farnoud and Milenkovic in [19]) *A minimum cost planar tree on a circle can be found in $O(r^4)$ time, where r is the number of vertices of a tree.*

Theorem 4 *φ_f -MCPS for a labeled circle on r vertices can be solved in $O(r^4)$ time.*

Conclusions and future directions

Practical matters

Our algorithm for φ_f -MCPS on a breakpoint graph with fixed labels has a running time of $O(n^5)$ in the worst case. Note that the running time is dominated, however, by the maximum bipartite matching step in “ φ -MCPS for a breakpoint graph” section. The size of the bipartite graph is determined by the number of AA and BB paths which is bounded by the maximum number of chromosomes m

for the two species. Thus using Lemma 4 we know that the algorithm scales like $O(mn^4)$ on biological data. For the same reason our algorithm for SORTING BY wDCJs AND INDELS IN INTERGENES [7] on a breakpoint graph scales like $O(m^2n \log n + m^3)$ instead of $O(n^3)$ on biological data. Further, n is the number of syntenic blocks—and not literally the genes as we call them. Our analyses of *Drosophila* genomes yield no AA and BB paths, and less than 100 blocks [10]. Our analysis of Human and Mouse genomes yields between 250 and 800 syntenic blocks, depending on the parameters given to OrthoCluster [28].

For graphs with higher degree nodes, like those graphs that represent genomes with duplicated genes, the number of simple cycles can grow rapidly. Although this relationship is beyond the scope of this work, we expect that fixed parameter algorithms could be developed to handle biological data in the future.

Future direction

Our cost framework is liberal, and in our examples we have explored only a small portion of its capacities. Edges can be labeled by more complex objects such as vectors or trees. The cost can be a function of a combination of the edge and vertex labels. We hope that a closer study of the graph $\mathcal{D}(G, \rho)$ from “Decomposition of a 2-break scenario” section will lead to polynomial time algorithms for φ -MCPS on circles for a large family of cost functions. Once the set of scenarios for a circle is better understood,

one could address the problems of counting and sampling the φ -MCPS scenarios.

While all of our results apply to genomes with insertions or deletions of single genes, further study is required in order to increase efficiency on genomes with duplicate genes.

Our assumption of “minimum evolution” may not always be true as an actual evolutionary scenario might be non-parsimonious [29]. The MINIMUM COST SCENARIO (MCS) problem of finding a minimum cost scenario among all the possible scenarios has already been studied for a couple of fairly simple cost functions [6, 12] and proven to be NP-hard in both of these cases. However, as we have shown in [12], computationally tractable algorithms can still be implemented for certain NP-hard MCS problems. An intermediate problem between MCPS and MCS could be the one of finding a minimum cost scenario among the scenarios of a prescribed length.

Proofs

Lemma 1

Lemma (Bienstock and Günlük in [20]) *The minimum length of a 2-break scenario on a graph G is $d_{2b}(G) = e(G) - c(G)$.*

Proof A 2-break can increase the size of a MAECD by at most 1 and the size of a MAECD of a terminal graph is $e(G)$. This leads to an inequality $d_{2b}(G) \geq e(G) - c(G)$.

In this paragraph the *length* of a cycle will be its number of black edges. For any cycle c of length $l > 1$ there is a 2-break transforming c into a union of length 1 and length $l - 1$ cycles. This way we obtain a scenario of length $l - 1$ for c , and can transform every cycle of a MAECD of G independently, obtaining a 2-break scenario of length $e(G) - c(G)$. Thus, $d_{2b}(G) \leq e(G) - c(G)$. \square

Lemma 2

Lemma (Yancopoulos et al. in [13]) *The minimum length of a DCJ scenario transforming genome A into B is equal to $d_{2b}(G(A, B)) = e(G(A, B)) - c(G(A, B))$.*

Proof $G(A, B)$ is constructed in such a way that for every DCJ $A \rightarrow A'$ the transformation $G(A, B) \rightarrow G(A', B)$ is a 2-break. Notably, a DCJ $\{a, b\} \rightarrow \{a\}, \{b\}$ results in a transformation $\{a, b\}, \{\circ, \circ\} \rightarrow \{a, \circ\}, \{b, \circ\}$, as the construction of a breakpoint graph guarantees that there are enough black loops $\{\circ, \circ\}$ to realize such a 2-break. For any 2-break $G(A, B) \rightarrow G'$ with $G' \neq G(A, B)$ there exists

a DCJ $A \rightarrow A'$ such that $G(A', B) = G'$. Since $G(B, B)$ is terminal, it follows that the minimum length of a scenario transforming A into B is $d_{2b}(G(A, B))$ and we conclude using Lemma 1. \square

Lemma 3

Lemma *If $\mathcal{D}(G, \rho)$ has k connected components then ρ can be partitioned into k subscenarios ρ^i and G can be partitioned into k edge-disjoint Eulerian subgraphs H^i in such a way that ρ^i is a scenario for H^i for every $i \in \{1, \dots, k\}$. If ρ is parsimonious, then $k = c(G)$ and $C(\rho) = \{H^1, \dots, H^k\}$ is a MAECD of G .*

Proof Take a connected component C of $\mathcal{D}(G, \rho)$. It has an equal number of vertices of indegree 0 and vertices of outdegree 0. Its edges incident to the vertices of indegree 0 are labeled with the black edges of G and its edges incident to the vertices of outdegree 0 are labeled with the gray edges of G . Together these labels define a subgraph H of G that we will prove to be Eulerian.

Define C_l to be a subgraph of $\mathcal{D}(G, \rho_l)$ consisting of its connected components containing the vertices of indegree 0 of C . This way $C_m = C$. Define H_l to be a subgraph of G_l containing the gray edges of H and the black edges of G_l labeling the edges of C_l incident to the vertices of outdegree 0. This way $H_0 = H$ and H_m is a terminal graph.

We prove that H is Eulerian by induction. H_m is Eulerian as it is terminal. Suppose that H_l is Eulerian. By construction the two edges of G_l replaced by the l -th 2-break of ρ either both belong to H_{l-1} or both are outside of H_{l-1} . In the first case, H_l is obtained from H_{l-1} via a 2-break and as H_l is Eulerian this means that H_{l-1} is also Eulerian. In the second case, $H_l = H_{l-1}$, thus the latter stays Eulerian. Thus $H = H_0$ is Eulerian and we obtain a subsequence of ρ that is a scenario for H .

$\mathcal{D}(G, \rho_0)$ has $e(G)$ connected components. The l -th 2-break of ρ merges two vertices of $\mathcal{D}(G, \rho_{l-1})$, thus reduces the number of the connected components by at most 1. This means that the number k of the connected components of $\mathcal{D}(G, \rho)$ is greater or equal to $e(G) - m$.

If ρ is parsimonious, then its length m is $e(G) - c(G)$ using Lemma 1. This means that $k \geq c(G)$ and G can be partitioned into k edge-disjoint Eulerian subgraphs. Due to the maximality of $c(G)$, we have that $k = c(G)$ and all of the obtained edge-disjoint Eulerian subgraphs of G are simple cycles. \square

Lemma 4

Lemma For a function f and an $O(f(r))$ time algorithm for φ -MCPS on a labeled circle on r vertices, there exists an $O(p^2f(P) + p^3 + f(n))$ time algorithm for φ -MCPS on a labeled breakpoint graph. If $f(r) = O(r^t)$ for some constant $t \geq 1$, then φ -MCPS on a labeled breakpoint graph can be solved in $O(pP^t + p^3 + n^t)$ time.

Proof The p^2 edges of a bipartite graph H can be weighted in $O(p^2f(P))$ time due to Theorem 2 and the fact that the simple cycles of $G(A, B)$ have at most 1 vertex of degree 2. A minimum weight maximum matching of H can be found in $O(p^3)$ time using the Hungarian algorithm. Finally, $MCPS_\varphi$ for the labeled circles in $G(A, B)$ can be computed in $O(f(n))$ time. Combining these results we obtain an $O(p^2f(P) + p^3 + f(n))$ time algorithm for computing $MCPS_\varphi(G(A, B), \lambda)$.

Now suppose that $f(r) = O(r^t)$ for some constant $t \geq 1$. Let a_1, \dots, a_p and b_1, \dots, b_p denote the number of edges in AA and BB paths with $\sum_{i=0}^p a_i = P_A, \sum_{j=0}^p b_j = P_B$ and $P = P_A + P_B$.

$MCPS_\varphi$ for a union of an AA path and a BB path having a and b edges respectively can be obtained by computing $MCPS_\varphi$ for at most two circles on $a + b$ vertices due to Theorem 2. This can be done in less than $c(a + b)^t$ steps for some constant c using the $O(r^t)$ time algorithm for computing $MCPS_\varphi$ for a circle. $MCPS_\varphi$ for every pair of AA and BB paths of $G(A, B)$ can be computed in a number of steps bounded by:

$$\begin{aligned} & \sum_{i=0}^p \sum_{j=0}^p c(a_i + b_j)^t \\ &= c \sum_{i=0}^p \sum_{j=0}^p \sum_{l=0}^t \binom{t}{l} a_i^l b_j^{t-l} c \sum_{l=0}^t \binom{t}{l} \sum_{i=For0}^p \sum_{j=0}^p a_i^l b_j^{t-l} \\ &= c \sum_{j=0}^p \sum_{i=0}^p b_j^t + c \sum_{i=0}^p \sum_{j=0}^p a_i^t + c \sum_{l=1}^{t-1} \binom{t}{l} \sum_{i=0}^p a_i^l \sum_{j=0}^p b_j^{t-l} \\ &= cp \sum_{j=0}^p b_j^t + cp \sum_{i=0}^p a_i^t + c \sum_{l=1}^{t-1} \binom{t}{l} \sum_{i=0}^p a_i^l \sum_{j=0}^p b_j^{t-l} \\ &\leq cp(\sum_{j=0}^p b_j)^t + cp(\sum_{i=0}^p a_i)^t + c \sum_{l=1}^{t-1} \binom{t}{l} (\sum_{i=0}^p a_i)^l (\sum_{j=0}^p b_j)^{t-l} \\ &\leq c(pP_B^t + pP_A^t) + pc \sum_{l=1}^{t-1} \binom{t}{l} P_B^{t-l} P_A^l \\ &= c(pP_B + pP_A)^t = cpP^t \end{aligned}$$

Thus, the weighting of H can be performed in $O(pP^t)$ time. This provides us with an $O(pP^t + p^3 + n^t)$ time algorithm for computing $MCPS_\varphi(G(A, B), \lambda)$. \square

Lemma 5

Lemma For a constant $t \geq 2$ and an $O(r^t)$ time α -approximation algorithm for φ -MCPS on a labeled circle on r vertices, there exists an $O(n^{t+1})$ time α -approximation algorithm for φ -MCPS on a labeled breakpoint graph.

Proof In Theorem 2, $MCPS_\varphi$ on a simple cycle is expressed as the minimum of the $MCPS_\varphi$ for a set of corresponding circles. In Theorem 1, $MCPS_\varphi$ on a graph is expressed as the minimum of the sums of the $MCPS_\varphi$ for the simple cycles. We prove an auxiliary proposition establishing the following:

1. An α -approximation for $MCPS_\varphi$ on a simple cycle can be obtained by taking the minimum of the α -approximations for the corresponding circles.
2. An α -approximation for $MCPS_\varphi$ on a graph can be obtained by taking the minimum of the sums of the α -approximations for $MCPS_\varphi$ on the simple cycles.

Proposition Take $k \in \mathbb{N}$ and two sets of positive numbers $\{q_1^*, \dots, q_k^*\}$ and $\{q_1, \dots, q_k\}$ with $q_i \leq \alpha q_i^*$ for every i . The following inequalities hold:

1. $\min\{q_i | i \in \{1, \dots, k\}\} \leq \alpha \min\{q_i^* | i \in \{1, \dots, k\}\}$
2. $\sum_{i=0}^k q_i \leq \alpha \sum_{i=0}^k q_i^*$

Proof Take u and v such that $q_u^* = \min\{q_i^* | i \in \{1, \dots, k\}\}$ and $q_v = \min\{q_i | i \in \{1, \dots, k\}\}$. By construction $q_v \leq q_u \leq \alpha q_u^*$ which proves the first inequality. For the second inequality it suffice to observe that $\sum_{i=0}^k q_i \leq \sum_{i=0}^k \alpha q_i^* = \alpha \sum_{i=0}^k q_i^*$. \square

A simple cycle of a breakpoint graph has at most one vertex of degree 2. This means that it has at most two corresponding circles (see Theorem 6). Taking the minimum of the α -approximations for $MCPS_\varphi$ on these circles provides us with an α -approximation for the simple cycle due to Theorem 6 and the first part of the proposition above. This way we obtain an α -approximation algorithm for φ -MCPS on a simple cycle of a breakpoint graph that

runs in $O(r^t)$ time where r is the number of the vertices in the simple cycle.

We can reuse the structure of a bipartite graph H presented in “ φ -MCPS for a breakpoint graph” section with the weights of the edges now being the α -approximations for the $MCPS_\varphi$ on the corresponding simple cycles. Following the same reasoning as in “ φ -MCPS for a breakpoint graph” section, we know that the minimum cost maximum matching of H leads to a MAECD of a breakpoint graph minimizing the sum of the α -approximations for the $MCPS_\varphi$ on its simple cycles. Combining Theorem 1, both parts of the proposition above, and the proof of Lemma 4, we obtain an $O(n^{t+1})$ time α -approximation algorithm for φ -MCPS on a breakpoint graph. \square

Lemma 6

Lemma *If ρ_O is a minimum 2-break-length O -scenario for a labeled circle (O, λ) , then $T(\rho_O)$ is a planar tree on (O, λ) . In addition to that, for a planar tree T on (O, λ) there exists an O -scenario ρ_O such that $T(\rho_O) = T$.*

Proof We prove the first statement by induction. It is trivially true if O has 2 vertices. We suppose it to be true for all the circles having less than $2l$ vertices and prove it for a circle having $2l$ vertices. Fix a minimum 2-break-length scenario ρ_O . Its length is $l - 1$ due to Lemma 1. The first O -break of ρ_O transforms (O, λ) into two vertex disjoint labeled circles (O_1, λ_1) and (O_2, λ_2) both having less vertices than O . The rest of the scenario ρ_O can be partitioned into ρ_O^1 acting on the edges of O_1 and ρ_O^2 acting on the edges of O_2 . As ρ_O is a minimum 2-break-length scenario, ρ_O^1 and ρ_O^2 must also be minimum 2-break-length scenarios. By the inductive hypothesis, $T(\rho_O^1)$ and $T(\rho_O^2)$ are planar trees on (O_1, λ_1) and (O_2, λ_2) respectively. $T(\rho_O)$ can be easily obtained from $T(\rho_O^1)$ and $T(\rho_O^2)$ by taking the union of their edges and adding an edge corresponding to the first 2-break of ρ_O . This way we obtain a planar tree $T(\rho_O)$ on (O, λ) proving the first statement of the lemma.

Now define the *distance* of an edge $\{x, y\}$ in T as the minimum number of vertices between x and y in the fixed circular embedding of T . For example, in the rightmost tree on the top of Fig. 5 the distance of the edge $\{w, z\}$ is one, because t is in between w and z , while the distance of the edge $\{x, y\}$ is 0. An edge is said to be *short* if its distance is 0. We prove an auxiliary proposition.

Proposition *A planar tree T on (O, λ) has a short edge incident to a leaf.*

Proof Choose a leaf x in T incident to an edge of the minimum distance d . If $d \neq 0$, then in between the leaf and the vertex that it is adjacent to, there are d other vertices. Since T is planar on (O, λ) , it is easy to see that there is at least one other leaf among these d vertices, which contradicts the minimality of x . \square

Now take a short edge $\{x, y\}$ incident to a leaf x in T . Take the black edges $\{u, v\}$ and $\{q, s\}$ in (O, λ) labeled with x and y respectively and separated by a gray edge $\{v, q\}$. Perform an O -break $(\{v, u\}, x), (\{q, s\}, y) \rightarrow (\{v, q\}, x), (\{u, s\}, y)$. resulting in two labeled circles. One of them is a terminal graph having two edges $\{v, q\}$ with the black one labeled with x . Remove the edge $\{x, y\}$ from T . This way we have reduced the size of the problem. The number of the vertices in the circle was reduced by two and the number of the edges in the tree was reduced by 1. We iterate this procedure to construct a required scenario. See the bottom part of Fig. 5 for an example. \square

Lemma 7

Lemma (Farnoud and Milenkovic in [19]) *A minimum cost planar tree on a circle can be found in $O(r^4)$ time, where r is the number of vertices of a tree.*

Proof Farnoud and Milenkovic pose the problem of sorting permutations by cost-constrained mathematical transpositions (a sorting scenario is called a *decomposition*) [19]. They define a cost function on the set of transpositions and treat the problem, called MIN-COST-MLD, of finding a minimum cost decomposition among the minimum length transposition decompositions of a permutation. They reduce this problem to finding a minimum cost planar tree on a circle, and propose the following $O(r^4)$ time dynamic programming algorithm for a tree having r vertices.

Enumerate the vertices 1 to r while respecting their order on the circle. Define $cost(i, j)$ as the minimum cost of a planar tree on the vertices $\{i, \dots, j\}$ for $1 \leq i < j \leq r$ and set $cost(i, i) = 0$ for $1 \leq i \leq r$.

Take a planar tree T on the vertices $\{1, \dots, r\}$. If $deg(1) = 1$ and 1 is on the edge $\{1, q\}$, then the cost of T is equal to $\Phi(1, q)$ plus the costs of the subgraphs of T induced by the vertices $\{2, \dots, q\}$ and $\{q + 1, \dots, r\}$. If $deg(1) > 1$, then take $q = \max(\{u | \{1, u\} \text{ belongs to } T\})$ and $s = \max(\{u | \text{there is a path in } T \text{ joining } 1 \text{ and } u \text{ but not visiting } q\})$. The cost of T is equal to $\Phi(1, q)$ plus the costs of the subgraphs of T induced by the vertices

$\{1, \dots, s\}$, $\{s + 1, \dots, q\}$ and $\{q, \dots, r\}$. This observation provides us with the following equality:

$$\text{cost}(i, j) = \max(\text{cost}(i, s) + \text{cost}(s + 1, q) + \text{cost}(q, j) + \Phi(i, q) \mid i \leq s < q \leq j)$$

for $1 \leq i < j \leq r$, that leads to an $O(r^4)$ time dynamic programming algorithm for finding $\text{cost}(1, r)$. \square

Acknowledgements

Not applicable.

Authors' contributions

All authors contributed to this work. All authors read and approved the final manuscript.

Funding

This work is partially supported by the IBC (Institut de Biologie Computationnelle) (ANR-11-BINF-0002), by the Labex NUMEV flagship project GEM, and by the CNRS project Osez l'Interdisciplinarité.

Availability of data and materials

Not applicable.

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Author details

¹ CNRS, LIRMM, Université Montpellier, 161 Rue Ada, 34392 Montpellier, France.

² Institut de Biologie Computationnelle (IBC), Montpellier, France.

Received: 29 April 2019 Accepted: 12 June 2019

Published online: 19 July 2019

References

- Blanchette M, Kunisawa T, Sankoff D. Parametric genome rearrangement. *Gene*. 1996;172(1):11–7.
- Baudet C, Dias U, Dias Z. Sorting by weighted inversions considering length and symmetry. *BMC Bioinform*. 2015;16(19):3.
- Biller P, Knibbe C, Guéguen L, Tannier E. Breaking good: accounting for the diversity of fragile regions for estimating rearrangement distances. *Genome Biol Evol*. 2016;8:1427–39.
- Nadeau JH, Taylor BA. Lengths of chromosomal segments conserved since divergence of man and mouse. *Proc Natl Acad Sci*. 1984;81(3):814–8.
- Ohno S. *Evolution by gene duplication*. Berlin: Springer; 1970. p. 160.
- Fertin G, Jean G, Tannier E. Algorithms for computing the double cut and join distance on both gene order and intergenic sizes. *Algorithms Mol Biol*. 2017;12(1):16.
- Bulteau L, Fertin G, Tannier E. Genome rearrangements with indels in intergenes restrict the scenario space. *BMC Bioinform*. 2016;17(14):426.
- Lieberman-Aiden E, Van Berkum NL, Williams L, Imakaev M, Ragoczy T, Telling A, Amit I, Lajoie BR, Sabo PJ, Dorschner MO, et al. Comprehensive mapping of long-range interactions reveals folding principles of the human genome. *Science*. 2009;326(5950):289–93.
- Veron A, Lemaitre C, Gautier C, Lacroix V, Sagot M-F. Close 3D proximity of evolutionary breakpoints argues for the notion of spatial synteny. *BMC Genomics*. 2011;12(1):303.
- Pulicani S, Simonaitis P, Rivals E, Swenson KM. Rearrangement scenarios guided by chromatin structure. In: RECOMB international workshop on comparative genomics. Berlin: Springer; 2017. p. 141–55.
- Swenson KM, Simonaitis P, Blanchette M. Models and algorithms for genome rearrangement with positional constraints. *Algorithms Mol Biol*. 2016;11(1):13.
- Simonaitis P, Swenson KM. Finding local genome rearrangements. *Algorithms Mol Biol*. 2018;13(1):9.
- Yancopoulos S, Attie O, Friedberg R. Efficient sorting of genomic permutations by translocation, inversion and block interchange. *Bioinformatics*. 2005;21(16):3340–6.
- Bergeron A, Mixtacki J, Stoye J. A unifying view of genome rearrangements. In: International workshop on algorithms in bioinformatics. Berlin: Springer; 2006. p. 163–73.
- Shao M, Lin Y. Approximating the edit distance for genomes with duplicate genes under DCJ, insertion and deletion. *BMC Bioinform*. 2012;13:13.
- Fosdick BK, Larremore DB, Nishimura J, Ugander J. Configuring random graph models with fixed degree sequences. *SIAM Rev*. 2018;60(2):315–55.
- Bafna V, Pevzner PA. Genome rearrangements and sorting by reversals. *SIAM J Comput*. 1996;25(2):272–89.
- Amir A, Levy A. String rearrangement metrics: a survey. *Algorithms and applications*. Berlin: Springer; 2010. p. 1–33.
- Farnoud F, Milenkovic O. Sorting of permutations by cost-constrained transpositions. *IEEE Trans Inf Theory*. 2012;58(1):3–23.
- Bienstock D, Günlük O. A degree sequence problem related to network design. *Networks*. 1994;24(4):195–205.
- Feder T, Guetz A, Mihail M, Saberi A. A local switch markov chain on given degree graphs with application in connectivity of peer-to-peer networks. In: FOCS'06. 47th annual IEEE symposium on foundations of computer science. 2006. p. 69–76.
- Caprara A. Sorting by reversals is difficult. In: Proceedings of the first annual international conference on computational molecular biology. New York: ACM; 1997. p. 75–83.
- Braga MDV, Sagot M-F, Scornavacca C, Tannier E. Bioinformatics research and applications: proceedings from ISBRA. The solution space of sorting by reversals. Berlin: Springer; 2007.
- Bitner JR. An asymptotically optimal algorithm for the dutch national flag problem. *SIAM J Comput*. 1982;11(2):243–62.
- Shao M, Lin Y, Moret BME. Sorting genomes with rearrangements and segmental duplications through trajectory graphs. *BMC Bioinform*. 2013;14:9.
- Birmelé E, Ferreira R, Grossi R, Marino A, Pisanti N, Rizzi R, Sacomoto G. Optimal listing of cycles and st-paths in undirected graphs. In: Proceedings of the twenty-fourth annual ACM-SIAM symposium on discrete algorithms. Society for Industrial and Applied Mathematics; 2013. p. 1884–96.
- Caprara A, Panconesi A, Rizzi R. Packing cycles and cuts in undirected graphs. In: European symposium on algorithms. Berlin: Springer; 2001. p. 512–23.
- Zeng X, Nesbitt MJ, Pei J, Wang K, Vergara IA, Chen N. Orthocluster: a new tool for mining synteny blocks and applications in comparative genomics. In: Proceedings of the 11th international conference on extending database technology: advances in database technology. New York: ACM; 2008. p. 656–67.
- Alexeev N, Alekseyev MA. Estimation of the true evolutionary distance under the fragile breakage model. *BMC Genomics*. 2017;18(4):356.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.