



**HAL**  
open science

# Assessing the robustness of parsimonious predictions for gene neighborhoods from reconciled phylogenies

Ashok Rajaraman, Cedric Chauve, Yann Ponty

## ► To cite this version:

Ashok Rajaraman, Cedric Chauve, Yann Ponty. Assessing the robustness of parsimonious predictions for gene neighborhoods from reconciled phylogenies. ISBRA - 11th International Symposium on Bioinformatics Research and Applications - 2015, Jun 2015, Norfolk, Virginia, United States. hal-01104587v1

**HAL Id: hal-01104587**

**<https://inria.hal.science/hal-01104587v1>**

Submitted on 18 Jan 2015 (v1), last revised 19 Mar 2015 (v3)

**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.

# Assessing the robustness of parsimonious predictions for gene neighborhoods from reconciled phylogenies

Ashok Rajaraman<sup>1,2</sup>, Cedric Chauve<sup>1</sup> and Yann Ponty<sup>1,3,4\*</sup>

<sup>1</sup>Department of Mathematics, Simon Fraser University, Burnaby (BC), Canada

<sup>2</sup>International Graduate Training Center in Mathematical Biology, Pacific Institute for Mathematical Sciences, Vancouver (BC), Canada

<sup>3</sup>CNRS/LIX, Ecole Polytechnique, Palaiseau, France

<sup>4</sup>Pacific Institute for the Mathematical Sciences, Vancouver (BC), Canada

Received on XXXXX; revised on XXXXX; accepted on XXXXX

Associate Editor: XXXXXXX

## ABSTRACT

**Motivations:** The availability of a large number of assembled genomes opens the way to study the evolution of syntenic character within a phylogenetic context. The DeCo algorithm, recently introduced by Bérard *et al.* allows the computation of parsimonious evolutionary scenarios for gene adjacencies, from pairs of reconciled gene trees. However, as for many phylogenetic algorithms, the robustness of a parsimonious scenario to a change in the cost associated to the considered evolutionary events is an important question, as different cost schemes might lead to significant different solutions. This is especially important for genome rearrangement models, that consider rare evolutionary events.

**Results:** Following the approach pioneered by Sturmfels and Pachter, we describe how to modify the DeCo dynamic programming algorithm into a parametric algorithm that allows the identification of classes of cost schemes that generates similar parsimonious evolutionary scenarios for gene adjacencies. Moreover, we extend this approach to assess the robustness (again to changes to the cost scheme of evolutionary events) of specific elements of evolutionary scenarios, here the presence/absence of specific ancestral gene adjacencies. We apply our method to a dataset of more than six thousands mammalian gene families, and show that computing the robustness to changes to cost schemes provides new and interesting insights on the evolution of gene adjacencies and the DeCo model.

**Availability:** <http://paleogenomics.irmacs.sfu.ca/DeClone/index.html>

**Contact:** [yann.ponty@lix.polytechnique.fr](mailto:yann.ponty@lix.polytechnique.fr)

## 1 INTRODUCTION

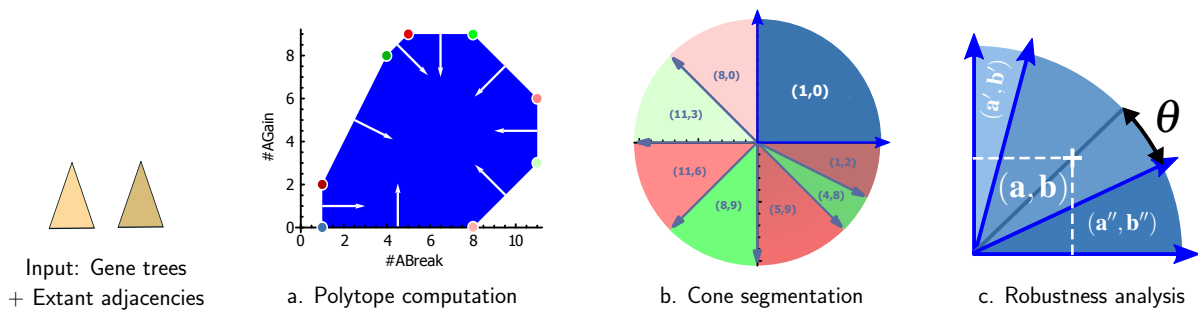
Reconstructing evolutionary histories of genomic characters along a given species phylogeny is a long-standing problem in computational biology. This problem has been studied for several types of genomic characters ranging from DNA sequences (genes and genomes segments) (Liberles, 2007) to gene family evolution (Bansal *et al.*, 2012; Doyon *et al.*, 2011) for which efficient algorithms exist to compute parsimonious evolutionary scenarios. Recently, Bérard *et al.* (2012) extended the corpus

of such results to syntenic characters. They introduced the notion of adjacency forests to model the evolution of gene adjacencies within a species phylogeny, together with an efficient dynamic programming (DP) algorithm, called DeCo, to compute parsimonious adjacency evolutionary histories. Reconstructing evolutionary scenarios for syntenic characters is an important step toward more comprehensive models of genome evolution, going beyond classical sequence/content frameworks as it implicitly integrates genome rearrangements (Chauve *et al.*, 2013). This may make available ancestral gene orders and genomes, whose usefulness in evolutionary biology has been demonstrated in several recent large-scale studies (Romanov *et al.*, 2014; Neafsey *et al.*, 2014)

The DeCo algorithm is a DP algorithm that takes as input a pair of reconciled gene trees, which represent the evolution of gene families in terms of speciation, duplication and gene loss within a given species tree, and a list of extant gene adjacencies, to compute, in polynomial time, an evolutionary scenario for adjacencies that minimizes the number of adjacency gains and breaks. The evolutionary model of DeCo is a combinatorial model that considers two events on adjacencies that can result from genome rearrangements: adjacency gain and adjacency break. To the best of our knowledge, DeCo is the only existing tractable model that considers the evolution of gene adjacencies within a phylogenetic framework accounting for both gene-specific events and genome rearrangement in adjacency evolution; so far other tractable models of genome rearrangements (eg. (Biller *et al.*, 2013; Tannier *et al.*, 2009)) accounting for a given species phylogeny are limited to single-copy genes and ignore gene-specific events.

The evolutionary events considered by DeCo, gene adjacency gain and break, caused by genome rearrangement, are rare evolutionary events, compared to gene-family specific events. For example, in a study of more than six thousand mammalian gene families, Bérard *et al.* used DeCo to reconstruct ancestral mammalian adjacencies, using a unit cost for both adjacency gain and adjacency break, and found an average of 1.25 adjacency gain/break per instance. It is then important to assess the robustness of inferences made by DeCo, whether it is of a parsimony cost, of an event count, or of an individual feature such as the presence of a specific ancestral adjacency. One way to address to this question is

\*to whom correspondence should be addressed



**Fig. 1.** Outline of our method for robustness analysis: Starting from two reconciled gene trees and a set of extant adjacencies, the polytope of parsimonious signatures is computed (a.). Its normal vectors define a segmentation of the space of cost schemes into cones (b.), each associated with a signature. In this example, further instantiated in Fig. 2, the positive quadrant is fully covered by a single cone, meaning that the parsimonious prediction does not depend on the precise cost model. In less extreme cases (c.), the robustness of a prediction (here, obtained using the  $(1, 1)$  scheme) to perturbations of the scheme can be measured as the smallest angle  $\theta$  such that a cost scheme at angular distance  $\theta$  no longer predicts a given signature  $(a, b)$ .

to consider the robustness of such features when compared to other possible evolutionary scenarios, an approach we explored recently by developing *DeClone*, a variant of *DeCo*, that considers the set of all possible evolutionary scenarios under a Boltzmann probability distribution. A second approach consists of considering changes in the cost associated to evolutionary events (the cost scheme) and in assessing how features of evolutionary scenarios are robust to such changes. Such approaches have recently been considered for the gene tree reconciliation problem and have been shown to improve significantly the results obtained from purely parsimonious approaches (Bansal *et al.*, 2013; Libeskind-Hadas *et al.*, 2014).

A second motivation is more specific to the *DeCo* model: as it models gene adjacencies, each ancestral gene can only be adjacent to at most two other genes, which is not considered in *DeCo*. However, the initial experiments using *DeCo* on mammalian gene trees resulted in a significant level of syntenic conflicts in the reconstructed ancestral adjacencies, as hundreds of ancestral genes were involved in more than two ancestral gene adjacencies. This raises the question of filtering inferred ancestral adjacencies to reduce the level of syntenic conflict, which can be done on the basis of their robustness. We reason that some of the erroneously-predicted adjacencies may result from an imperfectly calibrated evolutionary model. For instance, the cost scheme of *DeCo* consists of fixed penalties associated to the adjacency gain and break events, which are typically set to arbitrary integer values. Under this hypothesis, the features of a gene adjacency parsimonious evolutionary scenario that are not robust to perturbations of the cost scheme should be considered as dubious. This idea relates to the wider problem of deciding which precise cost to assign to evolutionary events in evolutionary models, a recurring question in the context of parsimony-based approaches in phylogeny.

These observations motivate the precise question tackled in this work: how robust is an inferred parsimonious ancestral gene adjacency to a change of the costs associated to adjacency gains and breaks? We address this problem using a methodology that can be traced back to parametric sequence alignment (Gusfield *et al.*, 1994), and has been formalized into a rigorous algebraic framework by Pachter and Sturmfels (Pachter and Sturmfels, 2004b,a). Since then, it has been applied to a few problems, notably of sequence alignment (Dewey *et al.*, 2006) and RNA folding

and alignment (Forouzmand and Chitsaz, 2013). As summarized in Fig. 1, in the context of the *DeCo* model, the main features of this approach, which we refer to as the *polytope approach* in the following, are (1) associating each evolutionary scenario to a *signature*, a vector of two integers  $(g, b)$  where  $g$  is the number of adjacency gains and  $b$  the number of adjacency breaks; and (2) partitioning the space of cost schemes (i.e. of the set of pairs of real non-negative numbers  $(x, y)$  where  $x$  is the cost associated to an adjacency gain and  $y$  the cost associated to an adjacency break) into convex regions such that, for all the cost schemes within a region, all optimal solutions obtained with such cost schemes have the same signature. This partition can be computed by an algorithm that is a direct translation of the DP algorithm into a polytope framework. We refer the interested reader to a book by Pachter and Sturmfels (2005) for a thorough description of the algebraic and algorithmic foundations of this general technique.

The contribution of the present paper is a polytope algorithm for the *DeCo* model. We also show how to estimate the support and robustness of specific ancestral gene adjacencies through a simple modification of the polytope algorithm. We apply our method to the same data set of over six thousand mammalian gene trees considered in (Bérard *et al.*, 2012) and show that using the robustness of parsimonious ancestral adjacencies to a change of the cost scheme provides valuable insight on the gene adjacency evolution problem.

## 2 PRELIMINARY: MODELS AND PROBLEMS

### 2.1 Models: from gene trees to adjacency forests

A *phylogeny* is a rooted tree which describes the evolutionary relationships of a set of elements (species, genes, ...) represented by its nodes: internal nodes correspond to ancestral elements, leaves to extant elements, and edges represent direct descents between parents and children. For a node  $v$  of a phylogeny, we denote by  $s(v)$  the species it belongs to. We consider here three kinds of phylogenies (illustrated in Figure 2): species trees, reconciled gene trees and adjacencies trees and forests.

The trees we consider are always rooted. For a tree  $T$  and a node  $x$  of  $T$ , we denote by  $T(x)$  the subtree rooted at  $x$ . If  $x$  is an internal node, we assume it has either one child, denoted by  $a(x)$ , or two

children, denoted by  $a(x)$  and  $b(x)$ . A tree where all internal nodes have two children is called a *binary tree*.

*Species tree and Reconciled gene trees.* A species tree  $S$  is a binary tree that describes the evolution of a set of related species from a common ancestor (the root of the tree), through the mechanism of *speciation*.

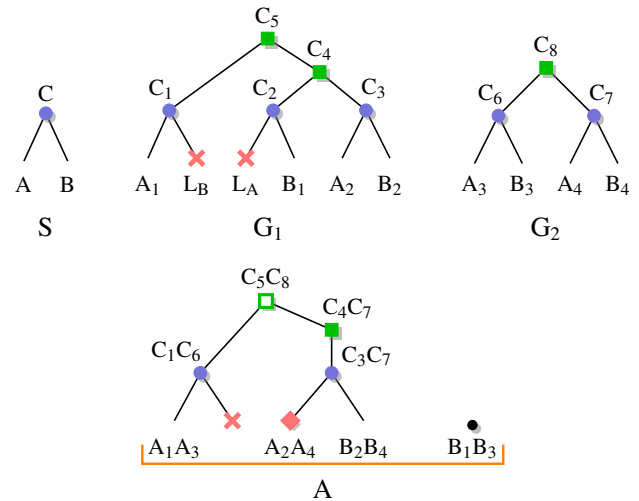
A reconciled gene tree is a binary tree that describes the evolution of a set of genes, called a *gene family*, through the evolutionary mechanisms of *speciation*, *gene duplication* and *gene loss*, within a given species tree  $S$ . Therefore, each node of a gene tree  $G$  represents either a gene loss, an extant gene or an ancestral gene. Ancestral genes are represented by the internal nodes of  $G$ , while gene losses and extant genes are represented by the leaves of  $G$ .

In a reconciled gene tree, we associate every ancestral gene (corresponding to an internal node  $g$ ) to an evolutionary event  $e(g)$  that leads to the creation of the two children  $a(g)$  and  $b(g)$ :  $e(g)$  is a *speciation* (denoted by **Spec**) if the species pair  $\{s(a(g)), s(b(g))\}$  is equal to the species pair  $\{a(s(g)), b(s(g))\}$ ,  $s(a(g)) \neq s(b(g))$ , or a *gene duplication* (**GDup**) if  $s(a(g)) = s(b(g)) = s(g)$ . If  $g$  is a leaf, then  $e(g)$  indicates either a *gene loss* (**GLoss**) or an *extant gene* (**Extant**), in which case  $e(g)$  is not an evolutionary event *stricto sensu*. A *pre-speciation* ancestral gene is an internal node  $g$  such that  $e(g) = \text{Spec}$ .

*Adjacency trees and forests.* We consider now that we are given two reconciled gene trees  $G_1$  and  $G_2$ , representing two gene families evolving within a species tree  $S$ . A *gene adjacency* is a pair of genes (one from  $G_1$  and one from  $G_2$ ) that appear consecutively along a chromosome, for a given species, ancestral or extant. Gene adjacencies evolve within a species tree  $S$  through the evolutionary events of *speciation*, *gene duplication*, *gene loss* (these three events, as described above, occur at the gene level and are modeled in the reconciled gene trees), and *adjacency duplication* (**ADup**), *adjacency loss* (**ALoss**) and *adjacency break* (**ABreak**), that are adjacency-specific events.

Following the model introduced in (Bérard *et al.*, 2012), we represent such an evolutionary history using an *adjacency forest*, composed of *adjacency trees*. An adjacency tree represents the evolution of an ancestral gene adjacency (located at the root of the tree) through the following events.

- The duplication of an adjacency  $\{g_1, g_2\}$ , where  $g_1$  and  $g_2$  are respectively genes from  $G_1$  and  $G_2$  such that  $s(g_1) = s(g_2)$ , follows from the simultaneous duplication of both its genes  $g_1$  and  $g_2$  (so  $e(g_1) = e(g_2) = \text{GDup}$ ), resulting in the creation of two distinct adjacencies each belonging to  $\{a(g_1), b(g_1)\} \times \{a(g_2), b(g_2)\}$ . This differs from the result of a single gene duplication (say gene  $e(g_1) = \text{GDup}$ ) that does not create a new adjacency but transforms  $\{g_1, g_2\}$  into either  $\{a(g_1), g_2\}$  or  $\{b(g_1), g_2\}$ .
- The loss of an adjacency, which can occur due to several events, such as the loss of exactly one of its genes (gene loss, **GLoss**), the loss of both its genes (adjacency loss, **ALoss**) or a genome rearrangement that breaks the contiguity between the two genes (adjacency break, **ABreak**).
- The creation/gain of an adjacency (denoted by **AGain**), for example due to a genome rearrangement, that results in the



**Fig. 2.** A species tree  $S$ , with two extant species  $A$  and  $B$  and an ancestral species  $C$ . Two reconciled gene trees  $G_1$  and  $G_2$ , with four extant genes in genome  $A$ , four extant genes in genome  $B$  and three ancestral genes in genome  $C$ . The set of extant gene adjacencies is  $(A_1A_3, B_1B_3, B_2B_4)$ . An adjacency forest  $A$  composed of two adjacency trees. Blue dots represent speciation nodes. Leaves are extant species/genes/adjacencies, except the one labeled by a red cross (gene loss) or a red diamond (adjacency breaks). Green squares are (gene or adjacency) duplication nodes. Gene labels refer to the species they belong to. Every node of the adjacency tree is labeled by a couple of nodes from gene trees representing a gene adjacency. Figure adapted from (Bérard *et al.*, 2012).

creation of a new adjacency tree whose root is the newly created adjacency.

*Evolution of gene families.* With this model, one can model the evolution of two gene families along a species phylogeny in terms of the evolutionary events introduced above by a triple  $(G_1, G_2, A)$ :  $G_1$  and  $G_2$  are reconciled gene trees representing the evolution of these families in terms of gene-specific events and  $A$  is an adjacency forest consistent with  $G_1$  and  $G_2$ . This notion of an evolutionary scenario is illustrated in Fig. 2.

Similar to species and reconciled gene trees, internal nodes of an adjacency tree are associated to ancestral adjacencies, while leaves are associated to extant adjacencies or lost adjacencies (due to a gene loss, adjacency loss or adjacency break), and are labeled by evolutionary events. The label  $e(v)$  of an internal node  $v$  of an adjacency forest  $A$  belongs to  $\{\text{Spec}, \text{GDup}, \text{ADup}\}$ , while the label  $e(v)$  of a leaf belongs to  $\{\text{Extant}, \text{GLoss}, \text{ALoss}, \text{ABreak}\}$ .

*Signatures, descriptors and parsimonious scenarios.* The signature (or event count) of an adjacency forest  $A$  is an ordered pair of integers  $\sigma(A) = (g_A, b_A)$  where  $g_A$  (resp.  $b_A$ ) is the number of adjacency gains (resp. adjacency breaks) in  $A$ . For example, the signature the adjacency forest displayed in Fig. 2 is  $(1, 1)$ . A *cost scheme* is a pair  $\mathbf{x} = (x_0, x_1)$  of non-negative real numbers, where  $x_0$  is the cost of an adjacency gain and  $x_1$  the cost of an adjacency break. The *cost* of an adjacency forest  $A$  for a given cost scheme  $\mathbf{x}$  is the number  $S(A) = x_0 \times g_A + x_1 \times b_A$ .

A *descriptor* of a scenario is a boolean or integer valued feature of the solution which does not contribute to the cost of the scenario, but rather represent a feature, a projection, of a scenario. For instance,

the presence/absence of an ancestral adjacency  $a = (g_1, g_2)$ , in a given adjacency forest  $A$  can be described as a boolean. Given  $k$  descriptors  $a_1, \dots, a_k$ , we can define an *extended signature* of a scenario  $A$  as a tuple  $\sigma_{a_1, \dots, a_k}(A) = (g, b, s_{a_1}, \dots, s_{a_k})$ , where  $g, b$  are the event counts for adjacency gains and breaks in  $A$  respectively, and  $s_{a_i}$  is the value of the descriptor  $a_i$  for  $A$ .

We consider now that we are given an evolutionary scenario  $(G_1, G_2, A)$  for two gene families as described above, together with a cost scheme  $\mathbf{x}$ . The adjacency forest  $A$  is *parsimonious* for  $\mathbf{x}$  if there is no other evolutionary scenario  $(G_1, G_2, B)$  such that  $S(B) < S(A)$ . The adjacency forest in Fig 2 is parsimonious for the cost scheme  $(1, 1)$ .

*Syntenic conflict.* Given a set of reconciled gene trees, together with a species tree, one can reconstruct ancestral gene orders by detecting all pairs of gene trees for which two extant genes form an extant gene adjacency – suggesting that there might have existed ancestral gene adjacencies with ancestral genes from these two families – and then compute a parsimonious adjacency forest for each such instance. Then ancestral gene orders can be obtained from ancestral adjacencies formed by pairs of pre-speciation ancestral genes.

However, if parsimonious adjacency forests are inferred independently, nothing prevents a pre-speciation ancestral gene from belonging to more than two ancestral adjacencies. This is obviously incompatible with the linear nature of the organization of genes along chromosomes and defines a *syntenic conflict*.

## 2.2 The DeCo and DeClone algorithms

Bérard *et al.* (2012) showed that, given a pair of reconciled gene trees  $G_1$  and  $G_2$ , a list  $L$  of extant gene adjacencies, and a cost scheme  $\mathbf{x}$ , one can compute, using a dynamic programming algorithm, an evolutionary scenario  $(G_1, G_2, A)$  where  $A$  is a parsimonious adjacency forest such that  $L$  is exactly the set of leaves of  $A$  labeled Extant.

The DeCo algorithm computes, for every pair of nodes  $g_1$  (from  $G_1$ ) and  $g_2$  (from  $G_2$ ) such that  $s(g_1) = s(g_2)$ , two quantities denoted by  $c_1(g_1, g_2)$  and  $c_0(g_1, g_2)$ , that correspond respectively to the cost of a parsimonious adjacency forest for the pairs of subtrees  $G(g_1)$  and  $G(g_2)$ , under the hypothesis that  $g_1$  and  $g_2$  do form (for  $c_1$ ) or do not form (for  $c_0$ ) an ancestral adjacency. As usual in dynamic programming along a species tree, the cost of a parsimonious adjacency forest for  $G_1$  and  $G_2$  is given by  $\min(c_1(r_1, r_2), c_0(r_1, r_2))$  where  $r_1$  is the root of  $G_1$  and  $r_2$  the root of  $G_2$ . Figure 1 in the appendix presents the DeCo dynamic programming equations.

In Chauve *et al.* (2014), we showed how to extend the DeCo algorithm to an ensemble approach that allows one to explore the whole solution space of adjacency forests. Let  $\mathcal{F}(G_1, G_2)$  be the set of all adjacency forests for  $G_1$  and  $G_2$ , including both optimal and sub-optimal ones. The *partition function* associated to  $G_1$  and  $G_2$  is defined by

$$\mathcal{Z}(G_1, G_2) = \sum_{A \in \mathcal{F}(G_1, G_2)} e^{-\frac{S(A)}{kT}}$$

where  $kT$  is an arbitrary constant. The partition function implicitly defines a *Boltzmann probability distribution* over  $\mathcal{F}(G_1, G_2)$ ,

where the probability of an adjacency forest  $A$  is defined by:

$$P(A) = \frac{e^{-\frac{S(A)}{kT}}}{\mathcal{Z}(G_1, G_2)}.$$

The Boltzmann probability of an adjacency, or more generally of a feature that can be observed in an adjacency forest, is then defined as the ratio of the sum of the Boltzmann probabilities of the adjacency forests that contain this feature. Such probabilities can be computed efficiently using a variation of the dynamic programming algorithm of DeCo (Chauve *et al.*, 2014). The impact of  $kT$  on the Boltzmann probability can be described as follows: when  $kT$  is small, the Boltzmann distribution probability is skewed toward parsimonious adjacency forests, while a high value of  $kT$  tends toward a more uniform probability distribution.

## 2.3 Robustness problems

The choice of a cost scheme is always a crucial step in evolutionary genomics. In the present work, we are generally interested in checking how a given feature of a parsimonious evolutionary scenario is robust to a change in the cost scheme. We define a cost scheme as a vector  $\mathbf{x} = (x_0, x_1) \in \mathbb{R}^2$  of positive real-valued numbers, each associated with a class of evolutionary events. Assuming that a cost scheme provides sufficient information to evaluate the objective function (here, the cost of an adjacency tree), then the predictions under such a model remain unchanged upon multiplying  $\mathbf{x}$  by any positive number. We may therefore assume that  $\|\mathbf{x}\| = 1$  without loss of generality. In two dimensions,  $\mathbf{x} = (x_0, x_1)$  can be summarized as a simple angle  $\theta$  (expressed in radians), and the difference between two cost schemes is indicated by their associated angular distance.

The first problem we are interested in is the *signature robustness problem*. We say that a signature  $\sigma = (g, b)$  is parsimonious for the cost scheme  $\mathbf{x}$  if there exists at least one adjacency forest  $A$  that is parsimonious for  $\mathbf{x}$  and has signature  $\sigma(A) = \sigma$ . Then, the robustness of the signature  $\sigma$  can be defined as the difference (i.e angular distance) from  $\mathbf{x}$  to the closest cost scheme for which  $\sigma$  is no longer parsimonious. The robustness of signatures can be computed using a modified version of the dynamic programming scheme of DeCo within an algebraic geometry framework adapted from (Pachter and Sturmfels, 2005), and described in further detail in the method section.

However, signatures only provide a quantitative summary of the evolutionary events described by a parsimonious adjacency forest. Similar signatures may be misleading, suggesting similar solutions while in fact corresponding to very dissimilar parsimonious adjacency forests. In particular, signatures discard any information about predicted sets of ancestral adjacencies. This leads us to consider the *parsimonious adjacency robustness problem*, where we address the robustness of inferred parsimonious ancestral adjacencies. Let  $a = (g_1, g_2)$  be an ancestral adjacency featured in a parsimonious adjacency forest for a cost scheme  $\mathbf{x}$ . We say that  $a$  is parsimonious for a cost scheme  $\mathbf{y}$  if  $a$  belongs to every adjacency forest that is parsimonious for  $\mathbf{y}$ . The robustness of  $a$  can then be defined as the angular distance from  $\mathbf{x}$  to the closest cost scheme  $\mathbf{y}$  for which  $a$  is no longer parsimonious. This notion of robustness can be extended to sets of ancestral adjacencies  $\mathbf{a}$  by considering that a set of ancestral adjacencies is parsimonious for a cost scheme  $\mathbf{y}$  if every adjacency in  $\mathbf{a}$  is parsimonious for  $\mathbf{y}$ .

### 3 METHODS

We now show how to compute the robustness measures introduced in the previous section. If the signature for a given adjacency forest  $A$  is given by the vector  $\sigma(A) = (g, b)$ , and the cost scheme is given by the vector  $\mathbf{x} = (x_0, x_1)$ , then the parsimony cost of  $\text{DeCo}$  can be written as the inner product  $\langle \mathbf{x}, \sigma(A) \rangle = g \times x_0 + b \times x_1$ .  $\text{DeCo}$  computes the following quantity for a pair of gene trees  $G_1$  and  $G_2$ .

$$c(G_1, G_2) = \min_{A \in \mathcal{F}(G_1, G_2)} \langle \mathbf{x}, \sigma(A) \rangle, \quad (1)$$

where  $\mathcal{F}(G_1, G_2)$  denotes the set all possible adjacency forests that can be constructed from  $G_1$  and  $G_2$ , irrespective of the cost scheme.

We will consider a single descriptor  $a$ , indicating the presence or absence of an ancestral adjacency  $a = (g_1, g_2) \in G_1 \times G_2$  of interest in the given instance, where  $s_a = 1$  if it is present in  $A$ , and 0 otherwise. Since, by definition, a descriptor does not contribute to the cost, we will only consider cost schemes of the form  $\mathbf{x} = (x_0, x_1, 0)$ , and  $\text{DeClone}$  will compute Eq. (1) as usual. However, the descriptor will allow us to compute the validity domains of  $a$ , as will be shown in next section.

At this point, it is important to note that more than one adjacency forest in the set  $\mathcal{F}(G_1, G_2)$  may be associated with the same signature, as signatures provide only an event count. Furthermore, for a given cost scheme  $\mathbf{x}$ , two adjacency forests  $A_1$  and  $A_2$  such that  $\sigma(A_1) = \sigma(A_2)$  have the same associated cost. We can thus define an equivalence class in  $\mathcal{F}(G_1, G_2)$  based on the signatures. However, the adjacency forests in this equivalence class may have different extended signatures, differing only in the last coordinate. Thus, there may be two adjacency forests  $A_1$  and  $A_2$  with extended signatures  $(g, b, 1)$  and  $(g, b, 0)$  respectively, and they will have the same cost for all cost schemes. Evolutionary scenarios with the same extended signature also form an equivalence class in  $\mathcal{F}(G_1, G_2)$ .

#### 3.1 Convex polytopes from signatures.

Let us denote the set of signatures of all scenarios in  $\mathcal{F}(G_1, G_2)$  by  $\sigma(\mathcal{F}(G_1, G_2))$ , and the set of extended signatures for a given adjacency  $a$  by  $\sigma_a(\mathcal{F}(G_1, G_2))$ . Each of these is a point in  $\mathbb{R}^d$ , where  $d = 2$  for signatures and  $d = 3$  for extended signatures. In order to explore the parameter space of parsimonious solutions to  $\text{DeCo}$ , we use these sets of points to construct a *convex polytope* in  $\mathbb{R}^d$ . A convex polytope is simply the set of all convex combinations of points in a given set, in this case the set of signatures or extended signatures (Pachter and Sturmfels, 2005). Thus, for each pair of gene trees  $G_1, G_2$  and a list of extant adjacencies, we can theoretically construct a convex polytope in  $\mathbb{R}^2$  by taking the convex combinations of all signatures in  $\sigma(\mathcal{F}(G_1, G_2))$ . This definition generalizes to a convex polytope in  $\mathbb{R}^3$  when extended signatures  $\sigma_a(\mathcal{F}(G_1, G_2))$  are considered for some signature  $a$ .

Viewing the set of evolutionary scenarios as a polytope allows us to deduce some useful properties:

1. Any (resp. extended) signature that is parsimonious for some cost scheme  $\mathbf{x}$  lies on the surface of the polytope;
2. If a (resp. extended) signature is parsimonious for two cost schemes  $\mathbf{x}$  and  $\mathbf{x}'$ , then it is also parsimonious for any cost scheme *in between* (for any convex combination of  $\mathbf{x}$  and  $\mathbf{x}'$ ).

Traditionally, a polytope is either represented as a set of inequations, which is inappropriate for our intended application. Therefore, we adopt a slightly modified representation, and denote the polytope of  $\mathcal{F}(G_1, G_2)$  as a list of signatures that are represented within  $\mathcal{F}(G_1, G_2)$  and lie on its convex hull.

A *vertex* in a polytope is a signature (resp. extended signature) which is the optimal for some cost scheme. The domain of parsimony of a vertex  $\mathbf{v}$  is the set of cost schemes for which  $\mathbf{v}$  is optimal. A consequence of Property 2 is that the domain of parsimony for a vertex  $\mathbf{v}$  is a *cone* in  $\mathbb{R}^d$ , formally defined as:

$$\text{Cone}(\mathbf{v}) = \left\{ \mathbf{x} \in \mathbb{R}^d : \langle \mathbf{x}, \mathbf{v} \rangle \leq \langle \mathbf{x}, \mathbf{w} \rangle \forall \mathbf{w} \in P \right\}. \quad (2)$$

The set of cones associated with the vertices of a polytope form a partition of the cost schemes space (Pachter and Sturmfels, 2005). This allows us to study the parameter space, and assess the effect of perturbing the parameters on the optimal solution of  $\text{DeCo}$ .

#### 3.2 Computing the polytope.

Pachter and Sturmfels (Pachter and Sturmfels, 2004b,a, 2005) introduced the concept of *polytope propagation*, based on the observation that the polytope of a dynamic programming (minimization) scheme can be computed through an algebraic substitution. Accordingly, any point that lies strictly within the polytope is suboptimal for any cost scheme, and can be safely discarded by a procedure that repeatedly computes the *convex hull*  $H(P)$  of the (intermediates) polytopes produced by the modified DP scheme. In the context of the  $\text{DeCo}$  DP scheme, the precise modifications, summarized in Table 1, are:

1. Any occurrence of the  $+$  operator is replaced by  $\oplus$ , the (convex) *Minkowski sum* operator, defined for  $P_1, P_2$  two polytopes as

$$P_1 \oplus P_2 = H(\{p_1 + p_2 \mid (p_1, p_2) \in P_1 \times P_2\});$$

2. Any occurrence of the  $\min$  operator is replaced by  $\cup$ , the *convex union* operator, defined for  $P_1, P_2$  two polytopes as

$$P_1 \cup P_2 = H(P_1 \cup P_2);$$

3. Any occurrence of an *adjacency gain* cost is replaced by the vector  $(1, 0)$  (resp.  $(1, 0, 0)$  for extended signatures);
4. Any occurrence of an *adjacency break* cost is replaced by the vector  $(0, 1)$  (resp.  $(0, 1, 0)$  for extended signatures);
5. (Extended signatures only) Any production that corresponds to the prediction of a fixed ancestral adjacency  $a$  in a scenario is replaced by the vector  $(0, 0, 1)$ ;

By making this substitution, we can efficiently compute the polytope associated with two input gene trees  $G_1$  and  $G_2$ , having sizes  $n_1$  and  $n_2$  respectively, through  $O(n_1 \times n_2)$  executions of the convex hull procedure. Let us remind that, in place of the integers used by the original minimization approach, intermediate convex polytopes are now processed by individual operations, and stored in the dynamic programming tables. The overall time and space complexities of the algorithm critically depend on the size of the polytopes, i.e. its number of vertices. Pachter and Sturmfels showed

DeCo algebra	Polytope propagation using signatures	Polytope propagation using extended signatures
Gain cost	(1, 0)	(1, 0, 0)
Break cost	(0, 1)	(0, 1, 0)
Adjacency presence	–	(0, 1, 0)
+	Minkowski Sum	Minkowski Sum
min	Convex hull	Convex hull
0	(0, 0)	(0, 0, 0)
$\infty$	Empty polytope	Empty polytope

**Table 1.** The operations in DeCo, with their analogous operations when moving to polytope propagation.

the important result that, in general, the number of vertices on the surface of the polytope is  $O(n^{d-1})$ , where  $d$  is the number of dimensions, and  $n$  is the size of the dynamic programming table. Thus, in our case, the number of vertices in the 2D polytope associated with simple signatures is in  $O(n_1 \times n_2)$ . This upper bound also holds for extended signatures, as the third coordinate is boolean, and the resulting 3D polytope is in fact the union of two 2D polytopes. The total cost of computing the convex hull is therefore bounded by  $O(n_1^2 \times n_2^2 \times \log(n_1 \times n_2))$ , e.g. using Chan’s convex hull algorithm (Chan, 1996).

As for the computation of the cones, let us remark that the cone of a vertex  $v$  in a given polytope  $P$  is fully delimited by a set of vectors, which can be computed from  $P$  as the normal vectors, pointing towards the center of mass of  $P$ , of each of the facets in which  $v$  appears. This computation can be performed as a postprocessing using simple linear algebra, and its consumption will remain largely dominated by that of the DP-fueled polytope consumption.

### 3.3 Assessing signature and adjacency robustness.

The cones associated with the polytopes cover all the real-valued cost schemes, including those associating negative costs to events. These cost schemes are not valid, and so, we only consider cones which contain at least one positive cost scheme.

*Signature robustness.* Given a fixed cost scheme  $\mathbf{y}$ , the vertex associated to the cone containing this cost scheme corresponds to the signature of all parsimonious scenarios for this cost scheme. In order to assess the robustness of this signature, we can calculate the smallest angular perturbation needed to move from  $\mathbf{y}$  to a cost scheme whose parsimonious scenarios do not have this signature. This is simply the angular distance from  $\mathbf{y}$  to the nearest boundary of the cone which contains it. Using this methods, we assign a numerical value to the robustness of the signatures of parsimonious scenarios on a number of instances for a particular cost scheme.

*Adjacency robustness.* In the case of extended signatures  $\sigma_a(\mathcal{F}(G_1, G_2))$  for an adjacency  $a$ , the polytope output by the DP equation is 3-dimensional. The cones associated with the vertices, as defined algebraically, now partition  $\mathbb{R}^3$ , the set of cost schemes  $(x_0, x_1, x_2)$ , where  $x_2$  indicates the cost of a distinguished adjacency. Since the third coordinate is a descriptor, it should not contribute to the cost scheme, and we must therefore restrict our analysis to the  $\mathbb{R}^+ \times \mathbb{R}^+ \times \{0\}$  subset of the cost scheme space. Therefore we take the intersection with the plane  $x_2 = 0$  of each

cone associated with a vertex  $(g, b, s_a)$ , and obtain the region in which the extended signature  $(g, b, s_a)$  is parsimonious. Note that this region is now a 2D cone.

However, the cost of an extended signature is now independent of the entry in its last coordinate, and there may exist two different extended signatures  $(g, b, 0)$  and  $(g, b, 1)$ , both parsimonious for all the cost schemes found in the 2D cone. It is also possible for adjacent cones to have different signatures, yet feature a given adjacency. The robustness of a given adjacency  $a$  is therefore computed from the cones using a greedy algorithm which, starting from the cone of  $\mathbf{x}$ , explores the adjacent cones in both directions (clockwise/counter-clockwise) until it finds one that no longer predicts  $a$ , i.e. is associated with at least one signature  $(g', b', 0)$ .

## 4 RESULTS AND DISCUSSION

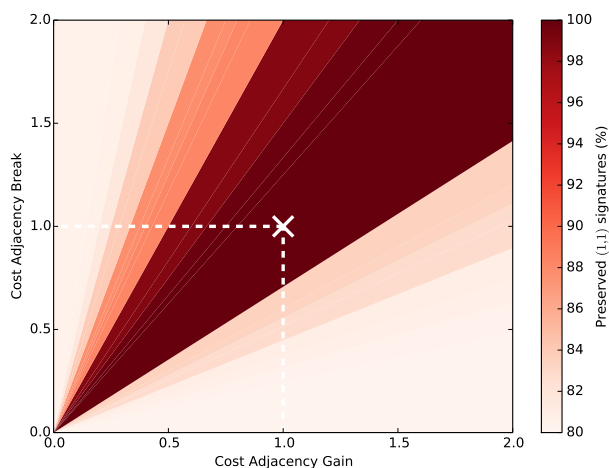
### 4.1 Data and DeCo analysis

The data set we considered is composed of 5,039 reconciled gene trees and 50,389 extant gene adjacencies, forming 6,074 DeCo instances, with genes taken from 36 extant mammalian genomes from the Ensembl database in 2012. In (Bérard *et al.*, 2012), these data were analyzed using the DeCo algorithm that computed a single parsimonious adjacency forest per instance. All together, these adjacency forests defined 96,482 ancestral adjacencies, covering 112,188 ancestral genes, where an “ancestral adjacency” is an adjacency that involves two genes  $g_1$  and  $g_2$  whose descendants in their respective gene trees do not belong to the same species  $s(g_1)$  (equal to  $s(g_2)$ ), i.e.  $g_1$  and  $g_2$  are pre-speciation genes, that were not duplicated within their species.

Besides our notions of robustness, an indirect validation criterion used to assess the quality of an adjacency forest is the limited presence of syntenic conflicts. An ancestral gene is said to participate to a syntenic conflict if it belongs to three or more ancestral adjacencies, as a gene can only be adjacent to at most two neighboring genes along a chromosome. An ancestral adjacency participates to a syntenic conflict if it contains a gene that does. Among the ancestral adjacencies inferred by DeCo, 16,039 participate to a syntenic conflict, covering 5,817 ancestral genes. This represents a significant level of syntenic conflict.

### 4.2 Results

We first considered all 6,074 instances, and computed for each signature the robustness of the parsimonious signature obtained with the cost scheme (1, 1) used in the DeCo experiment. Interestingly, we observe that for more than half of the instances, the parsimonious signature is robust to a change of cost scheme, as the cone associated to this signature is the complete first quadrant of the real plane. On the contrary, for 945 instances the parsimonious signature for the cost scheme (1, 1) is not robust to any change in the cost scheme; these cases correspond to interesting instances where the cost scheme (1, 1) lies at the border of two cones, meaning that two parsimonious signatures exist for the cost scheme (1, 1), and any small change of cost scheme tips the balance towards one of these two signatures. More generally, as revealed by Figure 3, the robustness landscape of the considered instances indicates an extreme robustness of parsimonious signatures. There is a  $\sim 80\%$



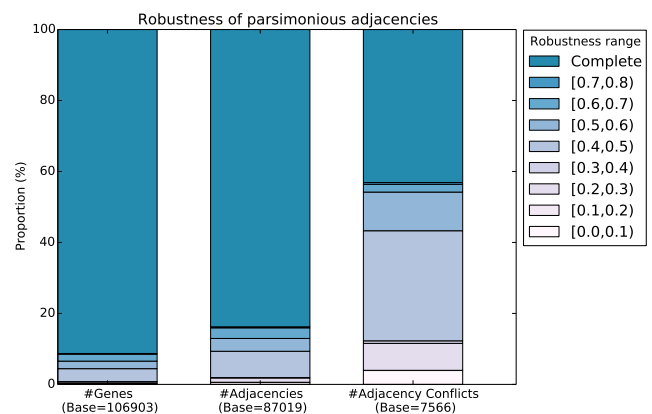
**Fig. 3.** Average robustness of signatures predicted using the  $(1, 1)$  cost scheme. At each point  $(x, y)$ , the colour indicates the proportion of signatures that are parsimonious, and therefore predicted, for the  $(1, 1)$  cost scheme, and remain parsimonious for the  $(x, y)$  cost scheme.

overlap between the sets of signatures that are parsimonious for any (positive) cost scheme, and for the  $(1, 1)$  cost scheme.

Finally, to evaluate the stability of the total number of evolutionary events inferred by parsimonious adjacency forests, we recorded three counts of evolutionary events for each instance: the total number of gene duplications in the reconciled trees, the number of syntenic events (adjacencies gains and breaks) of the parsimonious signature (called the parsimonious syntenic events count), and the maximum number of syntenic events taken over all signatures that are parsimonious for some cost scheme (called the maximum syntenic events count). We observe that the average number of gene duplications per instance is 3.38, and the average parsimonious (resp. maximum) syntenic events count is 1.25 (resp. 1.66). This shows a strong robustness of the number of syntenic events to changes in the cost scheme, with a significant difference in the number of syntenic events compared to gene-specific events.

We then considered the robustness of individual adjacencies. Using DeClone, we extracted adjacencies that belong to all parsimonious solutions for the cost scheme  $(1, 1)$  from all instances, and computed their robustness as defined in the previous sections. This set of ancestral adjacencies contains 87,019 adjacencies covering 106,903 ancestral genes, and participating in 7,566 syntenic conflicts. It can be observed that selecting such universally parsimonious ancestral adjacencies significantly reduced the number of syntenic conflicts, as almost all discarded ancestral adjacencies participated in syntenic conflicts.

The robustness of these adjacencies is summarized in Figure 4. It is interesting to observe that few adjacencies have a low robustness, while, conversely, a large majority of the universally parsimonious adjacencies are completely robust to a change of cost scheme (97,593 out of 106,639). This suggests that the DeCo model of parsimonious adjacency forests is robust, and infers highly supported ancestral adjacencies, which is reasonable given the relative sparsity of genome rearrangements in evolution compared to smaller scale evolutionary events. Considering syntenic conflicts,



**Fig. 4.** Universally parsimonious adjacencies and syntenic conflicts. (Left) Percentage of ancestral genes present in universally parsimonious adjacencies per level of minimum robustness of the adjacencies, expressed in radians. (Center) Percentage of universally parsimonious adjacencies per level of minimum robustness. (Right) Percentage of conserved conflicting ancestral adjacencies per level of minimum robustness.

we can notice a positive result, i.e. that filtering by robustness results in a significant decrease of the ratio of conflicting adjacencies. However it can be observed that even with robust universally parsimonious ancestral adjacencies, one can observe a significant number of adjacencies participating in syntenic conflicts. We discuss these observations in the next section.

### 4.3 Discussion

Our work raises several issues, both on methodological and applied aspects, that we discuss now.

From an application point of view, the ability to exhaustively explore the parameter space allowed to observe that, on the considered instances, the DeCo model is extremely stable. Even taking parsimonious signatures that maximize the number of evolutionary syntenic events (i.e. considering cost schemes that leads to the maximum number of events) results in an average increase of roughly 33% events (1.25 to 1.66), and stays very low, much lower than gene specific events such as gene duplications. This is consistent with the fact that for rare evolutionary events such as genome rearrangements, a parsimony approach is very relevant, especially when it can be complemented by efficient algorithms to explore slightly sub-optimal solutions, such as DeClone, and to explore the parameter space, as described in the present work. In terms of direct applications of the method developed here and in (Chauve *et al.*, 2014), one could think about a gene-tree based reconstruction of ancestral gene orders based on considering the set of all ancestral adjacencies that are parsimonious for at least one cost scheme (that can be computed in polynomial time using the method described in the present paper), scored using a mixture of their Boltzmann probability (that can be computed efficiently using DeClone) and robustness to changes of the cost scheme. This set of ancestral adjacencies is likely to contain many syntenic conflicts, that could be cleared out independently and efficiently for each ancestral species using the algorithm of Manuch *et al.* (2012).



An interesting observation is that even the set of ancestral adjacencies that are universally-parsimonious and robust to changes in the cost scheme contains a significant number of adjacencies participating in syntenic conflict. We believe the fact that this set of ancestral adjacencies is quite large supports the parsimony model of DeCo for studying the evolution of syntenic characters, as discussed above. We conjecture that the main reason for syntenic conflicts is in the presence of a significant number of erroneous reconciled gene trees. This is supported by the observation that the ancestral species with the highest number of syntenic conflict are also species for which the reconciliation with the mammalian species tree resulted in a significantly larger number of genes than expected (data not shown). This points clearly to errors in either gene tree reconstruction and in the reconciliation with the mammalian species phylogeny, that tends to assign wrong gene duplications in some specific species, resulting an inflation of the number of genes, especially toward the more ancient species (Hahn, 2007). It would be interesting to apply the same method we introduced to corrected mammalian gene trees (see (Lafond *et al.*, 2014) and references there regarding gene tree correction).

From a methodological point of view, there exists another way to explore the parameter space of a dynamic programming phylogenetic algorithm. It consists in computing, rather than optimal signatures for classes of cost schemes, the *Pareto-front* of the input instance (Libeskind-Hadas *et al.*, 2014; Saule and Giegerich, 2014). A signature  $v$  is said to be *Pareto-optimal* if there is no other signature whose entries are equal or smaller than the corresponding entries in  $v$ , and is strictly smaller at at least one coordinate. The Pareto-front is the set of all Pareto-optimal signatures, and can be efficiently computed by dynamic programming (Schnattinger *et al.*, 2013; Saule and Giegerich, 2014; Libeskind-Hadas *et al.*, 2014). In particular, Libeskind-Hadas *et al.* (2014) show how, in the context of phylogenetic tree reconciliation, this approach can be extended to provide a support for specific evolutionary events. The Pareto-front differs from the approach we describe in the present work in several aspects. An advantage of the Pareto-front is that it is a notion irrespective of the type of cost function being used. This contrasts with the polytope propagation technique, which requires that the cost function be a linear combination of its terms. However, so far, the Pareto-approach has only been used to define a partition of the parameter space when the cost function is restricted to be linear/affine, and it remains to investigate the difference with the polytope approach in this case. It is also key to note that Pareto-optimal signatures might correspond to points that are inside the polytope defined by parsimonious signatures, i.e. signatures of solutions that are not parsimonious for any cost scheme, and the polytope approach effectively ignores them rather than computing them. This difference deserves further investigation, especially to characterize Pareto-optimal but non-parsimonious signatures.

Finally there are several natural extensions of the present work that can be considered. The first one would be to consider a model of reconciled gene trees that accounts for lateral gene transfer. DeCo has been extended along this line in the case of a dated species phylogeny Patterson *et al.* (2013), and there does not seem to be any technical issue with adapting the polytope approach to this algorithm. Similarly, applying the Boltzmann and polytope approaches to the problem of gene tree reconciliations does not raise technical difficulties, and would complement existing approaches

that consider slightly different models (Bansal *et al.*, 2013; Libeskind-Hadas *et al.*, 2014). Another avenue, seemingly much more difficult, would be to integrate reconciliation and adjacencies. In the present work, we take as input reconciled gene trees, while there is no rationale, other than computational tractability, to separate gene specific events, such as gene duplication and gene loss, and syntenic events. However, it is likely that the question of simultaneously inferring reconciled gene trees and adjacency forests that minimize the total number of evolutionary events is difficult.

## 5 CONCLUSION

In the present work, we consider the robustness of the predictions of the DeCo (Bérard *et al.*, 2012) to perturbations of the evolutionary cost model algorithm. To that purpose, we adapt an algebraic geometry framework for dynamic programming, pioneered by Pachter and Sturmfels (2005), and use it to fully partition the space of cost schemes. We define objective, efficiently computed, measures of robustness for features (ancestral adjacency, event counts) of evolutionary scenario. Together with the ensemble approach described in (Chauve *et al.*, 2014), this provides tools that extend and complement the parsimonious approach to study the evolution of syntenic characters (Bérard *et al.*, 2012). Our analysis of a set of mammalian reconciled gene trees provides useful insights on this dataset, especially by revealing that a large majority of instances are extremely robust to a change in cost scheme, as are a large number of inferred ancestral adjacencies.

*Funding:* Research supported by an NSERC Discovery Grant to C.C., an SFU Michael Stevenson scholarship to A.R.

## REFERENCES

- Bansal, M. S., Alm, E. J., and Kellis, M. (2012). Efficient algorithms for the reconciliation problem with gene duplication, horizontal transfer and loss. *Bioinformatics*, **28**(12), 283–291.
- Bansal, M. S., Alm, E. J., and Kellis, M. (2013). Reconciliation revisited: Handling multiple optima when reconciling with duplication, transfer, and loss. *Journal of Computational Biology*, **20**(10), 738–754.
- Bérard, S., Gallien, C., Boussau, B., Szöllösi, G. J., Daubin, V., and Tannier, E. (2012). Evolution of gene neighborhoods within reconciled phylogenies. *Bioinformatics*, **28**(18), 382–388.
- Biller, P., Feijão, P., and Meidanis, J. (2013). Rearrangement-based phylogeny using the single-cut-or-join operation. *IEEE/ACM Trans. Comput. Biology Bioinform.*, **10**(1), 122–134.
- Chan, T. M. (1996). Optimal output-sensitive convex hull algorithms in two and three dimensions. *Discrete & Computational Geometry*, **16**, 361–368.
- Chauve, C., El-Mabrouk, N., Gueguen, L., Semeria, M., and Tannier, E. (2013). Duplication, rearrangement and reconciliation: A follow-up 13 years later. In *Models and Algorithms for Genome Evolution*, pages 47–62. Springer.
- Chauve, C., Ponty, Y., and Zanetti, J. P. P. (2014). Evolution of genes neighborhood within reconciled phylogenies: An ensemble approach. In *BSB*, volume 8826 of *Lecture Notes in Computer Science*, pages 49–56. Springer.

- Dewey, C. N., Huggins, P., Woods, K., Sturmfels, B., and Pachter, L. (2006). Parametric alignment of *Drosophila* genomes. *PLoS Computational Biology*, **2**(6).
- Doyon, J., Ranwez, V., Daubin, V., and Berry, V. (2011). Models, algorithms and programs for phylogeny reconciliation. *Briefings in Bioinformatics*, **12**(5), 392–400.
- Forouzmand, E. and Chitsaz, H. (2013). The RNA newton polytope and learnability of energy parameters. *Bioinformatics*, **29**(13), 300–307.
- Gusfield, D., Balasubramanian, K., and Naor, D. (1994). Parametric optimization of sequence alignment. *Algorithmica*, **12**(4/5), 312–326.
- Hahn, M. W. (2007). Bias in phylogenetic tree reconciliation methods: implications for vertebrate genome evolution. *Genome Biology*, **8**, R141.
- Lafond, M., Chauve, C., Dondi, R., and El-Mabrouk, N. (2014). Polytope refinement for the correction of dubious duplications in gene trees. *Bioinformatics*, **30**(17), 519–526.
- Liberles, D. A., editor (2007). *Ancestral Sequence Reconstruction*. Oxford University Press.
- Libeskind-Hadas, R., Wu, Y., Bansal, M. S., and Kellis, M. (2014). Pareto-optimal phylogenetic tree reconciliation. *Bioinformatics*, **30**(12), 87–95.
- Manuch, J., Patterson, M., Wittler, R., Chauve, C., and Tannier, E. (2012). Linearization of ancestral multichromosomal genomes. *BMC Bioinformatics*, **13**(S-19), S11.
- Neafsey, D. E. *et al.* (2014). Highly evolvable malaria vectors: The genomes of 16 anopheles mosquitoes. *Science*.
- Pachter, L. and Sturmfels, B. (2004a). Parametric inference for biological sequence analysis. *Proc Natl Acad Sci U S A*, **101**(46), 16138–16143.
- Pachter, L. and Sturmfels, B. (2004b). Tropical geometry of statistical models. *Proc Natl Acad Sci U S A*, **101**(46), 16132–16137.
- Pachter, L. and Sturmfels, B., editors (2005). *Algebraic Statistics for Computational Biology*. Cambridge University Press.
- Patterson, M., Szöllösi, G. J., Daubin, V., and Tannier, E. (2013). Lateral gene transfer, rearrangement, reconciliation. *BMC Bioinformatics*, **14**(S-15), S4.
- Romanov, M. *et al.* (2014). Reconstruction of gross avian genome structure, organization and evolution suggests that the chicken lineage most closely resembles the dinosaur avian ancestor. *BMC Genomics*, **15**(1), 1060.
- Saule, C. and Giegerich, R. (2014). Observations on the feasibility of exact pareto optimization. In *Proceedings of the 1st Workshop on Computational Methods for Structural RNAs (CMSR 2014), Strasbourg, France, September 7, 2014.*, pages 43–56.
- Schnattinger, T., Schöning, U., and Kestler, H. A. (2013). Structural RNA alignment by multi-objective optimization. *Bioinformatics*, **29**(13), 1607–1613.
- Tannier, E., Zheng, C., and Sankoff, D. (2009). Multichromosomal median and halving problems under different genomic distances. *BMC Bioinformatics*, **10**.