

# On the Gene Family-free DCJ Distance and Similarity

Fábio V. Martinez

Faculdade de Computação, Universidade Federal de Mato Grosso do Sul, Brazil

Algorithms in Comparative Genomics – Winter 2018/19

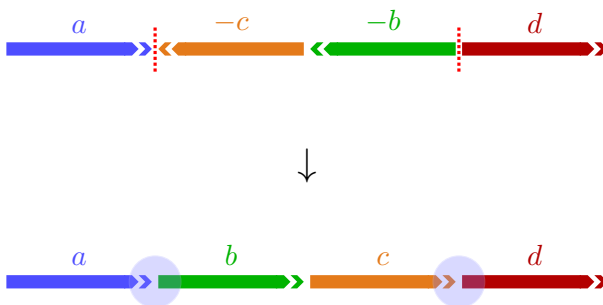
# Overview

- 1 Introduction
- 2 DCJ distance under the gene family-free method
- 3 DCJ similarity under the gene family-free method
- 4 Computational complexities
- 5 Algorithms
- 6 Experiments

# Genome rearrangements

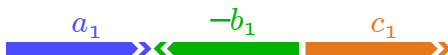
- ▶ Large-scale rearrangements change the number of chromosomes and/or the positions and orientations of genes (inversions, transpositions, translocations, TDRLs, fusions, fissions, ...)
- ▶ Genomes are represented as sequences of oriented DNA fragments (genes)

# The double-cut-and-join (DCJ) operation



# Problem

Classical problems: compute the DCJ distance/similarity between two given genomes



# Gene family-based approach

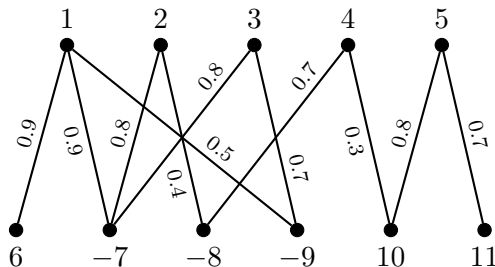
- ▶ each input genome contains one copy of each representative of a gene family (efficient algorithms)
- ▶ many copies of a gene present in input genomes (NP-hard problems)

# Gene family-free approach

- ▶ studying genome rearrangements without prior family assignment
- ▶ pairwise similarities between genes
- ▶ problems are not easier than problems under gene family-based approach

# The DCJ distance for the gene family-free method

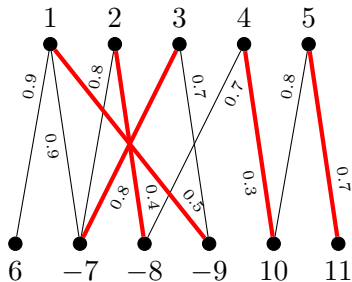
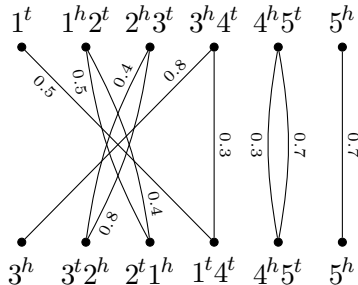
Gene similarity graph  $GS_\sigma(A, B)$



$\sigma : A \times B \rightarrow [0, 1]$  is the *normalized similarity function*



# The DCJ distance for the gene family-free method


 $GS_{\sigma}(A, B)$ 

 $AG_{\sigma}(A^M, B^M)$ 

$$d_{\sigma}(A^M, B^M) = 2|M| - \sigma(M) - c - i/2.$$

# The DCJ distance for the gene family-free method

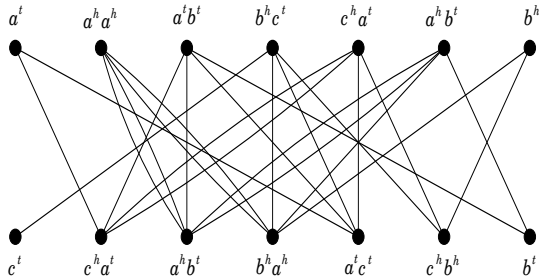
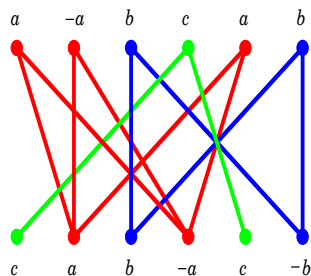
**Problem**  $\text{DIST}_{\text{DCJ}}^{\text{FF}}(A, B, \sigma)$ : Given genomes  $A$  and  $B$  and the gene similarity function  $\sigma$ , one can calculate the gene family-free DCJ distance between  $A$  and  $B$ :

$$d_{\text{DCJ}}^{\text{FF}}(A, B) = \min_{M \in \mathbb{M}} \{d_{\sigma}(A^M, B^M)\},$$

where  $\mathbb{M}$  is the set of all maximal matchings in  $GS_{\sigma}(A, B)$ .

# The DCJ distance for the gene family-based method

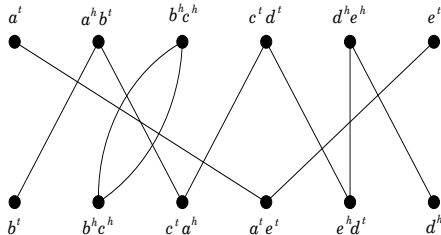
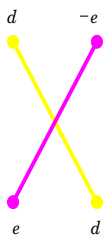
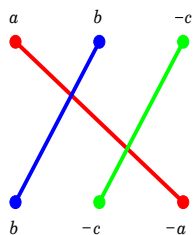
With duplicates



ILP-based algorithms, approximation algorithms

# The DCJ distance for the gene family-based method

Without duplicates

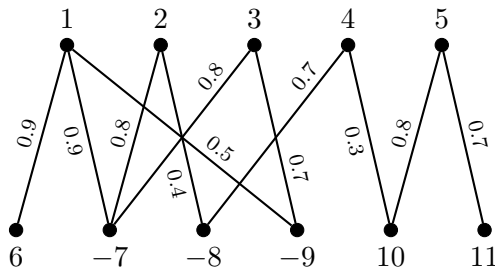


$AG(A, B)$

$$d_{\text{DCJ}}^{\text{FB}}(A, B) = n - c - i/2.$$

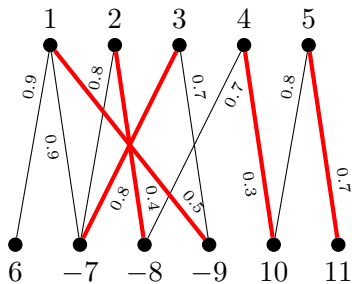
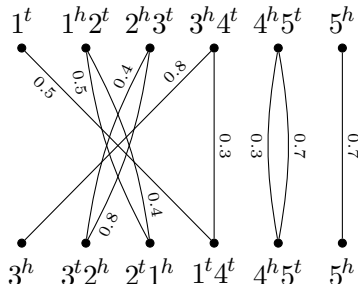
# The DCJ similarity for the gene family-free method

Gene similarity graph  $GS_\sigma(A, B)$



$\sigma : A \times B \rightarrow [0, 1]$  is the *normalized similarity function*

# The DCJ similarity for the gene family-free method


 $GS_{\sigma}(A, B)$ 

 $AG_{\sigma}(A^M, B^M)$

# The DCJ similarity for the gene family-free method

Let the *normalized weight*  $\widehat{w}(C)$  of a component  $C$  of  $AG_\sigma(A^M, B^M)$  be:

$$\widehat{w}(C) = \begin{cases} \frac{w(C)}{|C|}, & \text{if } C \text{ is a cycle,} \\ \frac{w(C)}{|C| + 1}, & \text{if } C \text{ is an odd path,} \\ \frac{w(C)}{|C| + 2}, & \text{if } C \text{ is an even path.} \end{cases}$$

Then, the wDCJ similarity is given by:

$$s_\sigma(A^M, B^M) = \sum_{C \in \mathcal{C}} \widehat{w}(C)$$

# The DCJ similarity for the gene family-free method

**Problem**  $\text{SIM}_{\text{DCJ}}^{\text{FF}}(A, B, \sigma)$ : Given genomes  $A$  and  $B$  and their gene similarities  $\sigma$ , calculate their gene family-free DCJ similarity

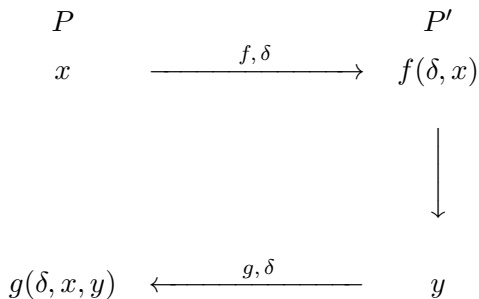
$$s_{\text{DCJ}}^{\text{FF}}(A, B) = \max_{M \in \mathbb{M}} \{s_{\sigma}(A^M, B^M)\},$$

where  $\mathbb{M}$  is the set of all maximal matchings in  $GS_{\sigma}(A, B)$ .



# Computational complexities

AP-reduction  $\leq_{AP}$

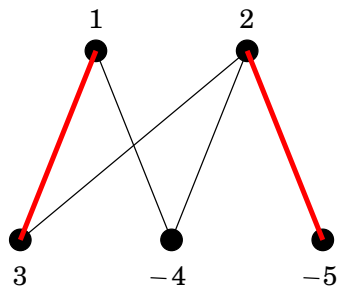


# Computational complexities

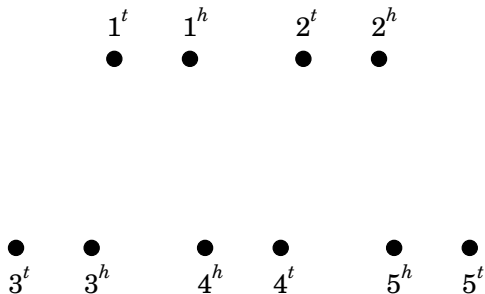
**Theorem**  $\text{DIST}_{\text{DCJ}}^{\text{FF}}$  is APX-hard and cannot be approximated with approximation ratio better than  $1237/1236 = 1.0008\dots$ , unless  $P = NP$ .

**Theorem**  $\text{SIM}_{\text{DCJ}}^{\text{FF}}$  is APX-hard and cannot be approximated with approximation ratio better than  $22/21 = 1.0476\dots$ , unless  $P = NP$ .

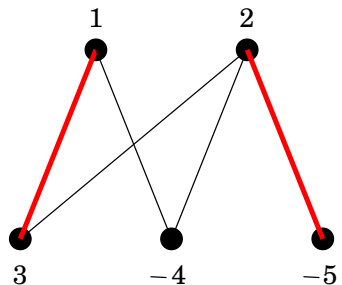
# Integer linear programs



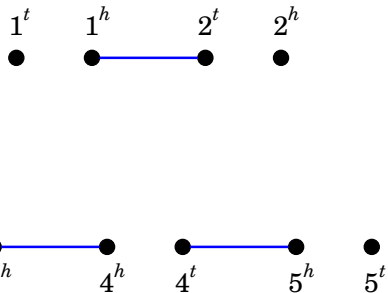
$$G = GS_{\sigma}(A, B) \quad M$$


 $H$

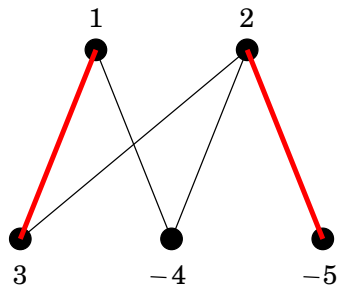
# Integer linear programs



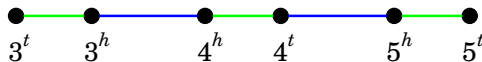
$$G = GS_{\sigma}(A, B) \quad M$$


 $H$

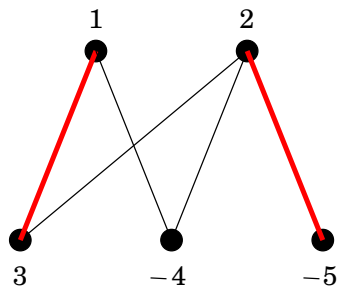
# Integer linear programs



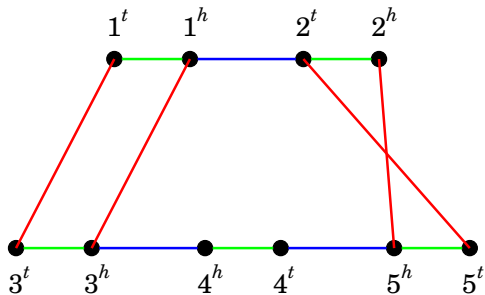
$$G = GS_{\sigma}(A, B) \quad M$$


 $H$

# Integer linear programs



$$G = GS_{\sigma}(A, B) \quad \mathbf{M}$$



$$H$$

# Integer linear program for $\text{DIST}_{\text{DCJ}}^{\text{FF}}$

$$\text{minimize} \quad 2 \sum_{e \in E_m} x_e - \sum_{e \in E_m} w_e x_e - \sum_{1 \leq i \leq |X_A|} z_i,$$

$$\begin{aligned} \text{subject to} \quad & x_e = 1 && \forall e \in E_a \\ & \sum_{uv \in E_m \cup E_s} x_{uv} = 1 && \forall u \in X_A \\ & \sum_{uv \in E_m \cup E_s} x_{uv} = 1 && \forall v \in X_B \\ & x_{ahbh} = x_{atbt} && \forall ab \in E(G) \\ & x_{ahat} + x_{bhbt} \leq 1 && \forall ab \in E(G) \\ & 0 \leq y_i \leq i && 1 \leq i \leq k \\ & y_i \leq y_j + i \cdot (1 - x_e) && \forall e = v_i v_j \in E(H) \\ & y_j \leq y_i + j \cdot (1 - x_e) && \forall e = v_i v_j \in E(H) \\ & i \cdot z_i \leq y_i && 1 \leq i \leq k \\ & i \cdot z_i \leq y_i && 1 \leq i \leq |X_A| \end{aligned}$$

# Integer linear program for $\text{DIST}_{\text{DCJ}}^{\text{FF}}$

$$\text{minimize} \quad 2 \sum_{e \in E_m} x_e - \sum_{e \in E_m} w_e x_e - \sum_{1 \leq i \leq |X_A|} z_i,$$

$$\begin{aligned} \text{subject to} \quad & x_e = 1 && \forall e \in E_a \\ & \sum_{uv \in E_m \cup E_s} x_{uv} = 1 && \forall u \in X_A \\ & \sum_{uv \in E_m \cup E_s} x_{uv} = 1 && \forall v \in X_B \\ & x_{ahbh} = x_{atbt} && \forall ab \in E(G) \\ & x_{ahat} + x_{bhbt} \leq 1 && \forall ab \in E(G) \\ & 0 \leq y_i \leq i && 1 \leq i \leq k \\ & y_i \leq y_j + i \cdot (1 - x_e) && \forall e = v_i v_j \in E(H) \\ & y_j \leq y_i + j \cdot (1 - x_e) && \forall e = v_i v_j \in E(H) \\ & i \cdot z_i \leq y_i && 1 \leq i \leq k \\ & i \cdot z_i \leq y_i && 1 \leq i \leq |X_A| \end{aligned}$$



# Integer linear program for $\text{DIST}_{\text{DCJ}}^{\text{FF}}$

$$\text{minimize} \quad 2 \sum_{e \in E_m} x_e - \sum_{e \in E_m} w_e x_e - \sum_{1 \leq i \leq |X_A|} z_i,$$

$$\begin{aligned} \text{subject to} \quad & x_e = 1 && \forall e \in E_a \\ & \sum_{uv \in E_m \cup E_s} x_{uv} = 1 && \forall u \in X_A \\ & \sum_{uv \in E_m \cup E_s} x_{uv} = 1 && \forall v \in X_B \\ & x_{ahbh} = x_{atbt} && \forall ab \in E(G) \\ & x_{ahat} + x_{bhbt} \leq 1 && \forall ab \in E(G) \\ & 0 \leq y_i \leq i && 1 \leq i \leq k \\ & y_i \leq y_j + i \cdot (1 - x_e) && \forall e = v_i v_j \in E(H) \\ & y_j \leq y_i + j \cdot (1 - x_e) && \forall e = v_i v_j \in E(H) \\ & i \cdot z_i \leq y_i && 1 \leq i \leq k \\ & i \cdot z_i \leq y_i && 1 \leq i \leq |X_A| \end{aligned}$$

# Algorithms for $\text{SIM}_{\text{DCJ}}^{\text{FF}}$

- ▶ ILP-based exact algorithm
- ▶ Heuristics
  - ▶ Maximum matching
  - ▶ Best density
  - ▶ Best length
  - ▶ Best length with weighted maximum independent set

# Experimental results for $\text{DIST}_{\text{DCJ}}^{\text{FF}}$

- ▶ data set of simulated genomes generated by ALF
- ▶ datasets with different genome sizes (1000, 2000 and 3000 genes) and evolutionary rates
- ▶ each dataset with 10 genomes, totalling 45 pairwise comparisons
- ▶ CPLEX was used to solve ILP instances with maximum running time was set to 60 minutes

# Experimental results for $\text{SIM}_{\text{DCJ}}^{\text{FF}}$

- ▶ data set of simulated genomes generated by ALF
  - ▶ datasets with different genome sizes (25, 50 and 1000 genes) and evolutionary rates
  - ▶ each dataset with 10 genomes, totalling 45 pairwise comparisons
  - ▶ Gurobi was used to solve ILP instances ILP with maximum running time was set to 60 minutes
  - ▶ Best density heuristic had the best performance
- ▶ real data set (human, house mouse and Norway rat)
  - ▶ 822, 953 and 863 genes, respectively
  - ▶ only heuristics
  - ▶ Best density heuristic had the best performance