# Combinatorial Optimization in Bioinfo Lecture 6 – Comparative genomics

Yann Ponty

AMIBio Team CNRS & École Polytechnique



Cool story bro. What next?

**Comparative genomics** studies (epi)genomics features of current (extant) species, to gain insight into their evolution and ancestral traits.

#### Some examples:

- Accurate inter-species distances by integrating genome-scale organization (genome rearrangements);
- Inference of ancestral features (genomic adjacencies);
- Reconciliations of species and gene trees...

Parsimony as a common framework

(but not exclusively!)

At **genomic scale** + **longer periods of time**, classic sequence homology becomes insufficient to infer elapsed time and precise correspondences.



[Pevzner and Tesler, Genome Research 2003]

How many inversions required to turn a Mouse into a Human? sort a signed permutation by reversals?

#### SORTING BY REVERSALS problem

**Input:** Signed permutation  $\pi$ 

**Output:** Shortest sequence of reversals turning  $\pi$  into 1, 2, 3 · · ·  $|\pi|$ 

Surprisingly, can be solved in polynomial-time [Hannenhalli–Pevzner, JACM 1999] Distance can even be obtained in linear time! [Baden-Moret-Yan, JCB 2004] At the origin of multiple works on sorting (structured/signed) permutations.

A simple(?) algorithm: [Hannenhalli–Pevzner 1999] + [Bergeron 2005]
 Mixed sign (+/-) π → Greedy reversals inducing max. #oriented pairs
 Positive sign π → Elimination (cut/merge → reversals) of hurdle(s)
 More details to come (lab assignment)

# Beyond (simple?) reversals



#### Fancy operations:

- ▶ Translocations: Exchange ends of chromosomes  $\rightarrow$  Linear time
- Franspositions: Exchange two consecutive sequences of genes  $\rightarrow$  NP-hard
- Double Cut and Join...
- Gain/loss of genetic material/syntenic blocks
  - ightarrow Handle multiple copies of genes
- Extension to unsigned permutations  $\rightarrow$  NP-hard.

Strong connection to graph theory and (enumerative) combinatorics

 $\rightarrow$  Very clever/elegant, but headache to extrapolate!



# Goal: To infer features of ancestral genomes



# Goal: To infer features of ancestral genomes



Phylogenetic tree T = (V, E) + Present/Absent **traits** for extant species What are the **most likely** features of ancestral species?

**Parsimony:** Find ancestors labeling  $\xi : V \to \{P, A\}$  minimizing total change:

$$f(\xi) = \sum_{(u,v)\in\mathcal{T}} d(\xi(u),\xi(v))$$



Phylogenetic tree T = (V, E) + Present/Absent **traits** for extant species What are the **most likely** features of ancestral species?

**Parsimony:** Find ancestors labeling  $\xi : V \to \{P, A\}$  minimizing total change:

$$f(\xi) = \sum_{(u,v)\in\mathcal{T}} d(\xi(u),\xi(v))$$



Phylogenetic tree T = (V, E) + Present/Absent **traits** for extant species What are the **most likely** features of ancestral species?

**Parsimony:** Find ancestors labeling  $\xi : V \to \{P, A\}$  minimizing total change:

$$f(\xi) = \sum_{(u,v)\in\mathcal{T}} d(\xi(u),\xi(v))$$



Phylogenetic tree T = (V, E) + Present/Absent **traits** for extant species What are the **most likely** features of ancestral species?

**Parsimony:** Find ancestors labeling  $\xi : V \to \{P, A\}$  minimizing total change:

$$f(\xi) = \sum_{(u,v)\in\mathcal{T}} d(\xi(u),\xi(v))$$



Phylogenetic tree T = (V, E) + Present/Absent **traits** for extant species What are the **most likely** features of ancestral species?

**Parsimony:** Find ancestors labeling  $\xi : V \to \{P, A\}$  minimizing total change:

$$f(\xi) = \sum_{(u,v)\in\mathcal{T}} d(\xi(u),\xi(v))$$

A labeling  $\xi$  is said to **extend**  $\xi'$  iff their labels coincide on extant species

#### PARSIMONY LABELING problem

**Input:** Species tree  $\mathcal{T}$ ; Labeling  $\xi' : V \to \{P, A\}$  of extant species **Output:** Complete labeling  $\xi^*$  extending  $\xi'$  which minimizes total change

$$\xi^{\star} = \underset{\xi \text{ ext. } \xi'}{\operatorname{argmin}} \sum_{(u,v) \in \mathcal{T}} d(\xi(u),\xi(v))$$

Solved by dynamic programming (hey, it's been a while!) [Sankoff/Rousseau 1975] based on a *hypothetical game*, *played* at any ancestral/internal node *u*:

- ▶ What if ancestral node *u* had a given label  $x \in \{P, A\}$ ?
- Then min change from u with ξ(u) = x is obtained by minimizing over possible labels y, y' for children v, v' of u, and summing:
  - Direct contribution to change: d(x, y) + d(x, y')
  - Min change of a labeling of subtrees v, v', independent of (u, x) (rec.)

$$\rightarrow f_x^{\star}(u) = \sum_{uv \in \mathcal{T}} \min_{y \in \{\mathsf{P},\mathsf{A}\}} d(x,y) + f_y^{\star}(v) \quad f_x^{\star}(u) = \begin{cases} 0 & \text{if } x = \xi'(u) \\ +\infty & \text{otherwise} \end{cases}$$



 $f_x^{\star}(u) = \sum_{uv \in \mathcal{T}} \min_{y \in \{P,A\}} d(x,y) + f_y^{\star}(v) \quad f_x^{\star}(u) = \begin{cases} 0 & \text{if } x = \xi'(u) \\ +\infty & \text{otherwise} \end{cases}$ 



 $f_x^{\star}(u) = \sum_{uv \in \mathcal{T}} \min_{y \in \{P,A\}} d(x,y) + f_y^{\star}(v) \quad f_x^{\star}(u) = \begin{cases} 0 & \text{if } x = \xi'(u) \\ +\infty & \text{otherwise} \end{cases}$ 



 $f_x^{\star}(u) = \sum_{uv \in \mathcal{T}} \min_{y \in \{P,A\}} d(x,y) + f_y^{\star}(v) \quad f_x^{\star}(u) = \begin{cases} 0 & \text{if } x = \xi'(u) \\ +\infty & \text{otherwise} \end{cases}$ 



 $f_x^{\star}(u) = \sum_{uv \in \mathcal{T}} \min_{y \in \{P, A\}} d(x, y) + f_y^{\star}(v) \quad f_x^{\star}(u) = \begin{cases} 0 & \text{if } x = \xi'(u) \\ +\infty & \text{otherwise} \end{cases}$ 



 $f_x^{\star}(u) = \sum_{uv \in \mathcal{T}} \min_{y \in \{P,A\}} d(x,y) + f_y^{\star}(v) \quad f_x^{\star}(u) = \begin{cases} 0 & \text{if } x = \xi'(u) \\ +\infty & \text{otherwise} \end{cases}$ 



**Toy model:** Min. #event  $\rightarrow d(\xi(u), \xi(v)) = 1$  if  $\xi(u) \neq \xi(v)$ ; 0 otherwise

 $f_{x}^{\star}(u) = \sum_{uv \in \mathcal{T}} \min_{y \in \{P,A\}} d(x,y) + f_{y}^{\star}(v) \quad f_{x}^{\star}(u) = \begin{cases} 0 & \text{if } x = \xi'(u) \\ +\infty & \text{otherwise} \end{cases}$ 

10/2<u>4</u>



**Toy model:** Min. #event  $\rightarrow d(\xi(u), \xi(v)) = 1$  if  $\xi(u) \neq \xi(v)$ ; 0 otherwise

 $f_{x}^{\star}(u) = \sum_{uv \in \mathcal{T}} \min_{y \in \{P,A\}} d(x,y) + f_{y}^{\star}(v) \quad f_{x}^{\star}(u) = \begin{cases} 0 & \text{if } x = \xi'(u) \\ +\infty & \text{otherwise} \end{cases}$ 

10/2<u>4</u>



**Toy model:** Min. #event  $\rightarrow d(\xi(u), \xi(v)) = 1$  if  $\xi(u) \neq \xi(v)$ ; 0 otherwise

 $f_{x}^{\star}(u) = \sum_{uv \in \mathcal{T}} \min_{y \in \{P,A\}} d(x,y) + f_{y}^{\star}(v) \quad f_{x}^{\star}(u) = \begin{cases} 0 & \text{if } x = \xi'(u) \\ +\infty & \text{otherwise} \end{cases}$ 

10/2<u>4</u>



**Toy model:** Min. #event  $\rightarrow d(\xi(u), \xi(v)) = 1$  if  $\xi(u) \neq \xi(v)$ ; 0 otherwise

 $f_{x}^{\star}(u) = \sum_{uv \in \mathcal{T}} \min_{y \in \{P,A\}} d(x,y) + f_{y}^{\star}(v) \quad f_{x}^{\star}(u) = \begin{cases} 0 & \text{if } x = \xi'(u) \\ +\infty & \text{otherwise} \end{cases}$ 



Sankoff/Rousseau parsimony represents a very general framework, which can be used to infer:

- presence/absence of regulatory motifs
- ancestral sequences
- ancestral structures (Sankoff algo. for RNA comparative folding)
- ancestral genome architectures
- ancestral allele frequencies

(adjacencies, gene loss/gain)

( $\infty$  #labels  $\rightarrow$  restriction on *d*)

Much can be learned from apparent contradictions...



How best to **explain** those topological discrepancies? and **learn** something along the way



Ancestral events (duplication/losses/...) induce apparent inconsistencies

**Reconciliations** represent evolutionary scenario that **explain** gene trees by embedding them within species trees **but** some more *reasonable* than others

Goal: Find parsimonious reconciliation



Ancestral events (duplication/losses/...) induce apparent inconsistencies

**Reconciliations** represent evolutionary scenario that **explain** gene trees by embedding them within species trees **but** some more *reasonable* than others

Goal: Find parsimonious reconciliation

# **Parsimonious DL reconciliations**

Duplication/Loss reconciliation: Binary tree R such that :

- leaf labeled with extant gene or Loss operation;
- internal node labeled with Dup or Spec event

All nodes are mapped through  $\sigma$  to a species

*R* congruent with gene tree *G* if removing Loss from *R* + contract 1<sup>ary</sup> branches  $\rightarrow G$ 

R congruent with species tree S if:

- For each Spec node  $u \to v$ , w in R, we have  $\sigma_u \to \sigma_v, \sigma_w$  in S
- For each Dup node u → v, w in R, we have σ<sub>u</sub> = σ<sub>v</sub> = σ<sub>w</sub>



#### DL RECONCILIATION problem

Input: Gene tree G, Species tree S

Output: Reconciliation R congruent to both G and S, minimizing cost

 $Cost(R) = \alpha \times \#Loss(R) + \beta \times \#Dup(R)$ 

Can be solved efficiently using LCA mapping... [Page, Syst. Biol. 1994] ... or using dynamic programming in  $\mathcal{O}(|S| \times |G|)$  time/space

DL RECONCILIATION **problem** 

Input: Gene tree G, Species tree S

Output: Reconciliation R, congruent to G and S, minimizing cost

 $\mathsf{Cost}(R) = \alpha \times \#\mathsf{Loss}(R) + \beta \times \#\mathsf{Dup}(R)$ 

# Notations:

C[g, s] Best cost of reconciliation for node/subtrees  $g \sqsubseteq G$  and  $s \sqsubseteq S$  $L_u, R_u$  Left/right child of node u

 $\sigma(u)$  Species of a leaf *u* within gene tree (or species tree...duh!)

$$C[g, s] = \begin{cases} 0 & \text{if } \sigma(g) = \sigma(s) \\ +\infty & \text{otherwise} \end{cases} \begin{array}{l} C[g, s] \\ g \mid \text{eaf,s internal} \end{array} = +\infty \\ \\ C[g, s] = \begin{cases} \beta + C[L_g, s] + C[R_g, s] & \text{Ancestor gene } g \text{ duplicated in species } s \\ C[L_g, L_s] + C[R_g, R_s] & \text{Speciation event, direct mapping} \\ C[L_g, R_s] + C[R_g, L_s] & \text{Speciation event, crossed mapping} \\ \alpha + C[g, L_s] & \text{Loss of gene subtree } g \text{ in right part of species tree} \\ \alpha + C[g, R_s] & \text{Loss of gene subtree } g \text{ in left part of species tree} \end{cases}$$

- Efficient algorithms exist for parsimonious scenarios for genomic characters (genes under DTL, for instance).
- Bérard et al (2012) extended methods to gene adjacencies under gene Speciation, Duplication, Loss.
- ▶ DeCo considers cost of gains/ breaks due to rearrangements.

**Problem:** Given two reconciliations  $(R_1, R_2)$  + extant adjacencies, construct parsimonious evolutionary scenario



#### Output: Adjacency forest



#### Event costs/parameters

- Adjacency gain: x
- Adjacency break: y



## Output: Adjacency forest



#### Event costs/parameters

- Adjacency gain: x
- Adjacency break: y



1.  $e(q_1) = \text{Extant}$  and  $e(q_2) = \text{Extant}$ :  $c_1(g_1, g_2) = \begin{cases} 0 & \text{if } g_1g_2 \text{ is an adjacency} \\ 0 & \text{otherwise} \end{cases}$   $c_0(g_1, g_2) = \begin{cases} 0 & \text{if } g_1g_2 \text{ is not an adjacency} \\ 0 & \text{otherwise} \end{cases}$  e(g<sub>1</sub>) = GLoss and e(g<sub>2</sub>) ∈ {Extant, Spec, GDup}:  $c_1(g_1, g_2) = c_0(g_1, g_2) = 0^{\#NonGDup(g_2)}$ 3.  $e(g_1) \in \{\text{Extant}, \text{Spec}, \text{GDup}\} \text{ and } e(g_2) = \text{GLoss}:$  $c_1(g_1, g_2) = c_0(g_1, g_2) = 0^{\# NonGDup(g_1)}$ 4.  $e(q_1) = \text{GLoss and } e(q_2) = \text{GLoss}$ :  $c_1(g_1, g_2) = c_0(g_1, g_2) = 0$ 5.  $e(q_1) \in \{\text{Extant}, \text{Spec}\}\ \text{and}\ e(q_2) = \text{GD}$ 
$$\begin{split} c_1(g_1,g_2) &= \min \begin{cases} c_1(g_1,b_{22}) + c_0(g_1,a_{22}) + x, \\ c_1(g_1,g_2) = g_1(g_1,g_2) + x, \\ c_1(g_1,g_2) = g_2(g_1,g_2) + x, \\ c_0(g_1,g_2) = g_2(g_1,g_2) + g_2(g_1,g_2) + x, \\ c_0(g_1,g_2) + g_2(g_1,g_2) + g_2(g_1,g_2) + x, \\ c_1(g_1,b_{22}) + c_0(g_1,a_{22}) + x, \\ c_1(g_1,b_{22}) + c_0(g_1,b_{22}) + x, \\ c_1(g_1,b_{22}) + c_0(g_1,b_{22}) + x, \\ c_1(g_1$$
6.  $e(q_1) = \mathsf{GDup}$  and  $e(q_2) \in \{\mathsf{Extant}, \mathsf{Spec}\}$ : 
$$\begin{split} c_1(g_1,g_2) &= \min \begin{cases} c_1(a_g_1,g_2) + c_0(b_g_1,g_2), & c_0(a_g_1,g_2) + c_1(b_g_1,g_2), \\ c_1(a_g_1,g_2) + c_1(b_g_1,g_2) + x, & c_0(a_g_1,g_2) + c_0(b_g_1,g_2) + y \end{cases} \\ c_0(g_1,g_2) &= \min \begin{cases} c_0(a_g_1,g_2) + c_0(b_{g_1},g_2), & c_0(a_g_1,g_2) + c_1(b_{g_1},g_2) + x, \\ c_1(a_g_1,g_2) + c_0(b_{g_1},g_2), & c_0(a_{g_1},g_2) + c_1(c_{g_1},g_2) + x, \end{cases} \end{split}$$
7.  $e(q_1) =$ Spec and  $e(q_2) =$ Spec:  $c_1(g_1,g_2) = \min \begin{cases} c_1(a_{g_1},b_{g_2}) + c_1(b_{g_1},a_{g_2}) + 0, & c_1(a_{g_1},b_{g_2}) + c_0(b_{g_1},a_{g_2}) + y + 0, \\ c_0(a_{g_1},a_{g_2}) + c_1(b_{g_1},a_{g_2}) + y + 0, & c_0(a_{g_1},b_{g_2}) + c_0(b_{g_1},a_{g_2}) + 2y + 0, \\ c_1(a_{g_1},a_{g_2}) + c_1(b_{g_1},b_{g_2}) + 0, & c_1(a_{g_1},a_{g_2}) + c_0(b_{g_1},b_{g_2}) + y + 0, \\ c_0(a_{g_1},a_{g_2}) + c_1(b_{g_1},b_{g_2}) + y + 0, & c_0(a_{g_1},a_{g_2}) + c_0(b_{g_1},b_{g_2}) + y + 0, \end{cases}$  $c_0(g_1,g_2) = \min \begin{cases} c_0(a_1,a_2) + c_1(y_0,y_2) + y \neq v, & c_0(y_0,y_1,y_2) + c_0(y_0,y_1,y_2), & ..., \\ c_0(a_2,b_2) + c_0(b_{21},a_{22}) + c_0(b_{21},a_{22}) + c_0(b_{21},b_{22}) + c_0(b_{21},a_{22}) + x + 0, \\ c_0(a_2,b_2) + c_1(b_{21},a_{22}) + x + 0, & c_1(a_1,b_2) + c_1(b_{21},a_{22}) + 2x + 0, \\ c_0(a_2,a_{22}) + c_0(b_{21},b_{22}) + x + 0, & c_1(a_2,b_{22}) + c_0(b_{21},b_{22}) + x + 0, \\ c_0(a_2,a_{22}) + c_0(b_{21},b_{22}) + x + 0, & c_1(a_2,a_{22}) + c_0(b_{21},b_{22}) + x + 0, \\ c_0(a_2,a_{22}) + c_0(b_{21},b_{22}) + x + 0, & c_1(a_2,a_{22}) + c_0(b_{21},b_{22}) + x + 0, \\ c_0(a_2,a_{22}) + c_0(b_{21},b_{22}) + x + 0, & c_1(a_2,a_{22}) + c_0(b_{21},b_{22}) + x + 0, \\ c_0(a_2,a_{22}) + c_0(b_{21},b_{22}) + c_0(b_{21},b_{22}) + c_0(b_{21},b_{22}) + x + 0, \\ c_0(a_2,a_{22}) + c_0(b_{21},b_{22}) + c_0(b_{21},b_{22}) + c_0(b_{21},b_{22}) + x + 0, \\ c_0(a_2,a_{22}) + c_0(b_{21},b_{22}) + c_0(b_{21},b_{22}) + c_0(b_{21},b_{22}) + x + 0, \\ c_0(a_2,a_{22}) + c_0(b_{21},b_{22}) + c_0(b_{21},b_{22})$ 

 $c_0(a_{g_1}, a_{g_2}) + c_1(b_{g_1}, b_{g_2}) + x + 0$ ,  $c_1(a_{g_1}, a_{g_2}) + c_1(b_{g_1}, b_{g_2}) + 2x + 0$ ,

8.  $e(q_1) = \mathsf{GDup}$  and  $e(q_2) = \mathsf{GDup}$ :

		$(c_1(a_{g_1}, g_2) + c_0(b_{g_1}, g_2)),$	$c_0(a_{g_1}, g_2) + c_1(b_{g_1}, g_2),$
	$c_1(g_1,g_2)=\min \phi$	$c_1(a_{g_1}, g_2) + c_1(b_{g_1}, g_2) + x,$	$c_0(a_{g_1}, g_2) + c_0(b_{g_1}, g_2) + y,$
		$c_1(g_1, a_{g_2}) + c_0(g_1, b_{g_2}),$	$c_0(g_1, a_{g_2}) + c_1(g_1, b_{g_2}),$
		$c_1(g_1, a_{g_2}) + c_1(g_1, b_{g_2}) + x,$	$c_0(g_1, a_{g_2}) + c_0(g_1, b_{g_2}) + y,$
		$c_1(a_{g_1}, a_{g_2}) + c_1(b_{g_1}, b_{g_2}) + c_0$	$(a_{g_1}, b_{g_2}) + c_0(b_{g_1}, a_{g_2}),$
		$c_1(a_{g_1}, a_{g_2}) + c_1(b_{g_1}, b_{g_2}) + c_0$	$(a_{g_1}, b_{g_2}) + c_1(b_{g_1}, a_{g_2}) + x,$
		$c_1(a_{g_1}, a_{g_2}) + c_1(b_{g_1}, b_{g_2}) + c_1$	$(a_{g_1}, b_{g_2}) + c_0(b_{g_1}, a_{g_2}) + x,$
		$c_1(a_{g_1}, a_{g_2}) + c_1(b_{g_1}, b_{g_2}) + c_1$	$(a_{g_1}, b_{g_2}) + c_1(b_{g_1}, a_{g_2}) + 2x,$
		$c_1(a_{g_1}, a_{g_2}) + c_0(b_{g_1}, b_{g_2}) + c_0$	$(a_{g_1}, b_{g_2}) + c_0(b_{g_1}, a_{g_2}) + y,$
		$c_1(a_{g_1}, a_{g_2}) + c_0(b_{g_1}, b_{g_2}) + c_0$	$(a_{g_1}, b_{g_2}) + c_1(b_{g_1}, a_{g_2}) + x + y$
		$c_1(a_{q_1}, a_{q_2}) + c_0(b_{q_1}, b_{q_2}) + c_1$	$(a_{q_1}, b_{q_2}) + c_0(b_{q_1}, a_{q_2}) + x + y$
		$c_0(a_{q_1}, a_{q_2}) + c_1(b_{q_1}, b_{q_2}) + c_0$	$(a_{q_1}, b_{q_2}) + c_0(b_{q_1}, a_{q_2}) + y,$
		$c_0(a_{q_1}, a_{q_2}) + c_1(b_{q_1}, b_{q_2}) + c_0$	$(a_{q_1}, b_{q_2}) + c_1(b_{q_1}, a_{q_2}) + x + y$
		$c_0(a_{q_1}, a_{q_2}) + c_1(b_{q_1}, b_{q_2}) + c_1$	$(a_{q_1}, b_{q_2}) + c_0(b_{q_1}, a_{q_2}) + x + y$
		$c_0(a_{q_1}, a_{q_2}) + c_0(b_{q_1}, b_{q_2}) + c_1$	$(a_{q_1}, b_{q_2}) + c_1(b_{q_1}, a_{q_2}),$
		$c_0(a_{q_1}, a_{q_2}) + c_1(b_{q_1}, b_{q_2}) + c_1$	$(a_{q_1}, b_{q_2}) + c_1(b_{q_1}, a_{q_2}) + x,$
		$c_1(a_{a_1}, a_{a_2}) + c_0(b_{a_1}, b_{a_2}) + c_1$	$(a_{a_1}, b_{a_2}) + c_1(b_{a_1}, a_{a_2}) + x,$
		$c_0(a_{a_1}, a_{a_2}) + c_0(b_{a_1}, b_{a_2}) + c_1$	$(a_{a_1}, b_{a_2}) + c_0(b_{a_1}, a_{a_2}) + y_i$
		$c_0(a_{a_1}, a_{a_2}) + c_0(b_{a_1}, b_{a_2}) + c_0$	$(a_{a_1}, b_{a_2}) + c_1(b_{a_1}, a_{a_2}) + y_1$
		$c_0(a_{a_1}, a_{a_2}) + c_0(b_{a_1}, b_{a_2}) + c_0$	$(a_{a_1}, b_{a_2}) + c_0(b_{a_1}, a_{a_2}) + 2y,$
	$c_0(g_1, g_2) = \min \{$	$(c_0(a_{a_1}, a_2) + c_0(b_{a_1}, a_2))$	$c_0(a_{a_1}, a_2) + c_1(b_{a_1}, a_2) + x_1$
		$c_1(a_{a_1}, a_2) + c_0(b_{a_1}, a_2) + x_1$	$c_1(a_{a_1}, a_2) + c_1(b_{a_1}, a_2) + 2x$
		$c_0(a_1, a_{a_2}) + c_0(a_1, b_{a_2})$	$c_0(a_1, a_{a_2}) + c_1(a_1, b_{a_2}) + x.$
		$c_1(a_1, a_2) + c_2(a_1, b_2) + \tau$	$c_1(a_1, a_2) + c_1(a_1, b_2) + 2r$
		$(v_1(g_1) \circ g_2) + o((g_1) \circ g_2) + o)$	~1(g1) ~g2/ · ~1(g1) ~g2/ + at.

Deco algorithm [Berard et al. 2012] **Dyn. prog. scheme** for most parsimonious adjacency forest in  $\mathcal{O}(|R_1||R_2|)$ .

$$\begin{array}{l} \textbf{7.} \text{If } e(g_1) = \text{Spec and } e(g_2) = \text{Spec:} \\ \textbf{7.} \text{If } e(g_1) = \text{Spec and } e(g_2) = \text{Spec:} \\ \textbf{7.} \text{If } e(g_1) = \text{Spec and } e(g_2) = \text{Spec:} \\ \textbf{7.} \text{If } e(g_1) = \text{Spec and } e(g_2) = \text{Spec:} \\ \textbf{7.} \text{If } e(g_1) = \text{Spec and } e(g_2) + c_1(b(g_1), a(g_2)), c_1(a(g_1), b(g_2)) + c_0(b(g_1), a(g_2)) + \text{ABreak}, \\ \textbf{7.} \text{C}_0(a(g_1), b(g_2)) + c_1(b(g_1), a(g_2)) + 2\text{ABreak}, \\ \textbf{7.} \text{C}_1(a(g_1), a(g_2)) + c_1(b(g_1), b(g_2)) + 2\text{ABreak}, \\ \textbf{7.} \text{C}_0(a(g_1), a(g_2)) + c_1(b(g_1), b(g_2)) + 2\text{ABreak}, \\ \textbf{7.} \text{C}_0(a(g_1), a(g_2)) + c_1(b(g_1), b(g_2)) + 2\text{ABreak}, \\ \textbf{7.} \text{C}_0(a(g_1), b(g_2)) + c_1(b(g_1), b(g_2)) + 2\text{ABreak}, \\ \textbf{7.} \text{C}_0(a(g_1), b(g_2)) + c_1(b(g_1), a(g_2)) + 2\text{ABreak}, \\ \textbf{7.} \text{C}_0(a(g_1), b(g_2)) + c_1(b(g_1), a(g_2)) + 2\text{ABreak}, \\ \textbf{7.} \text{C}_0(a(g_1), b(g_2)) + c_1(b(g_1), a(g_2)) + 2\text{ABreak}, \\ \textbf{7.} \text{C}_0(a(g_1), b(g_2)) + c_0(b(g_1), a(g_2)) + 2\text{ABreak}, \\ \textbf{7.} \text{C}_0(a(g_1), b(g_2)) + c_0(b(g_1), a(g_2)) + 2\text{ABreak}, \\ \textbf{7.} \text{C}_0(a(g_1), a(g_2)) + c_0(b(g_1), a(g_2)) + 2\text{ABreak}, \\ \textbf{7.} \text{C}_0(a(g_1), a(g_2)) + c_0(b(g_1), a(g_2)) + 2\text{ABreak}, \\ \textbf{7.} \text{C}_0(a(g_1), a(g_2)) + c_0(b(g_1), a(g_2)) + 2\text{ABreak}, \\ \textbf{7.} \text{C}_0(a(g_1), a(g_2)) + c_0(b(g_1), a(g_2)) + 2\text{ABreak}, \\ \textbf{7.} \text{C}_0(a(g_1), a(g_2)) + c_0(b(g_1), a(g_2)) + 2\text{ABreak}, \\ \textbf{7.} \text{C}_0(a(g_1), a(g_2)) + c_0(b(g_1), a(g_2)) + 2\text{ABreak}, \\ \textbf{7.} \text{C}_0(a(g_1), a(g_2)) + c_0(b(g_1), b(g_2)) + 2\text{ABreak}, \\ \textbf{7.} \text{C}_0(a(g_1), a(g_2)) + c_0(b(g_1), b(g_2)) + 2\text{ABreak}, \\ \textbf{7.} \text{C}_0(a(g_1), a(g_2)) + c_0(b(g_1), b(g_2)) + 2\text{ABreak}, \\ \textbf{7.} \text{C}_0(a(g_1), a(g_2)) + 2\text{ABreak}, \\ \textbf{7.} \text{C}_0(a(g_1$$

**Deco algorithm** [Berard *et al*, 2012] **Dyn. prog. scheme** for most parsimonious adjacency forest in  $O(|R_1| \times |R_2|)$ 

**Intuition:** Jointly explores both reconciliations in all possible ways, to account for adjacency breaks/gains.

## Comparative genomics is a vibrant field with strong 2-ways links with

- Phylogenetics, (Epi)genomics, Paleogenomics, Proteomics, Structural biology, Systems Biology ...
- ...but also string processing, enumerative combinatorics, discrete algorithms, discrete maths...

It predicts the past, but also informs studies of the unobservable present

Some  $\pm$  current directions:

- Beyond parsimony, Bayesian approaches (likelihood maximization)
   Mainly based on sampling (ABC), but also more direct opt. (DP)
- Distinguishing between erroneous and informative incongruences
- Expressive models for reconciliations: Horizontal gene transfers, polytomies...
- Integrating multiple steps of analyses into joint optimizations e.g. comparative genome scaffolding
- Advanced algorithms: FPT algo/kernels, guaranteed approximations, exp. time (small exponents), ILP... ... for all of the above (+ML)

# Onwards to our (final) lab assignment 1/4

**Goal:** Implement sorting by reversal for signed permutations [Bergeron 2006] w.l.o.g. assume that input permutation  $\pi$  flanked by 0 and n+1



Some definitions:

Oriented pair: Two pos. (i, j) with consecutive values of opposite signs



**Reversal** of oriented pairs (i, j):

$$\pi_{i} + \pi_{j} = +1 \rightarrow \cdots \qquad \pi_{i} \qquad \pi_{i+1} \qquad \cdots \qquad \pi_{j-1} \qquad \pi_{j} \qquad \cdots \\ \dots \qquad -\pi_{j-1} \qquad -\pi_{j-2} \qquad \cdots \qquad -\pi_{i} \qquad \pi_{j} \qquad \cdots \\ \pi_{i} + \pi_{j} = -1 \rightarrow \qquad \cdots \qquad \pi_{i} \qquad \pi_{i+1} \qquad \cdots \qquad \pi_{j-1} \qquad \pi_{j} \qquad \cdots \\ \dots \qquad \pi_{i} \qquad -\pi_{i} \qquad \cdots \qquad -\pi_{i+2} \qquad -\pi_{i+1} \qquad \cdots$$

Goal: Implement sorting by reversal for signed permutations [Bergeron 2006] w.l.o.g. assume that permutations flanked by 0 and n+1

Definitions for positive permutations:

Framed interval: Interval [i, j] in  $\pi$  of the form

 $\pi_i \pi_{i+1} \pi_{i+2} \cdots \pi_j$  with  $\pi_j = \pi_i + (j-i)$  and  $\pi_i < \pi_k < \pi_j, \forall i < k < j$ 

or **equivalently** region [i, j] contains a permutation of  $[\pi_i, \pi_j]$ 

But... one needs to consider circular ordering of values



► Hurdles: Minimal framed intervals w.r.t. inclusion

 0
 2
 5
 4
 3
 6
 1
 7

 Hurdles: [2, 6], [6, 2]

Goal: Implement sorting by reversal for signed permutations [Bergeron 2006] w.l.o.g. assume that permutations flanked by 0 and n+1

Definitions for positive permutations:

Splitting a hurdle [i, j]: Look for position k such that  $\pi_k = \pi_i + 1$ , and reverse interval [i + 1, k - 1]

Input: 
$$0 2 4 3 1 5$$
  
Single hurdle: [1, 6]  
 $0 2 4 3 1 5$   
Output:  $0 -3 -4 -2 1 5$ 

• Merging two hurdles [i, j] and  $[k, l], j \le k$ : Reverse interval [j, k]

Goal: Implement sorting by reversal for signed permutations [Bergeron 2006]

**Final algo.** While  $\pi$  not sorted:

- If  $\pi$  negative:
  - ▶ Identify Oriented Pair (OP) (*i*, *j*) that creates max #OPs
  - Reverse OP (i, j)
  - Iterate until permutation  $\pi$  becomes positive
- If  $\pi$  positive, consider hurdles  $\mathcal{H}$  of  $\pi$ :

$$\mathcal{H} = \{(i, j)\} \\ \rightarrow \text{split} (i, j)$$

- $\mathcal{H} = \{(i, j), (k, l)\}$  $\rightarrow \text{merge } (i, j)$
- $|\mathcal{H}| > 2, |\mathcal{H}| \text{ even}$ 
  - $\rightarrow$  Merge two non-consecutive hurdles

 $|\mathcal{H}| > 2, |\mathcal{H}| \text{ odd}$ 

**Simple hurdle**  $h^*$ : Cutting  $h^*$  + reversal of OPs **decreases** #hurdles

- If  $\mathcal H$  contains simple hurdle  $h^*$ 
  - $\rightarrow$  Cut  $h^*$
- Otherwise
  - $\rightarrow$  Merge two non-consecutive hurdles (or consecutive ones if  $|\mathcal{H}|=3)$