

SCJ with duplicated genes (distance, median, SPP)

Cedric Chauve

Department of Mathematics, Simon Fraser University

LaBRI, Université de Bordeaux

Bielefeld Block Seminar on Comparative Applied Genomics

From the previous lecture, we know that, in the case of genomes with no duplications,

- The SPP is tractable solely in the SCJ model, with the issue that it considers only a restricted solution space.
- The weighted SPP is tractable only when $\alpha = 0$ (MWM) or $\alpha = 1$ (pure SCJ). FPT algorithm based on Dynamic Programming (DP).
- In the SCJ-TD-FD distance, the directed median is tractable, the rooted median is NP-hard.

In this lecture, we will explore in details what can be done in the SCJ model toward usable solutions of genome rearrangement problems.

Our first step toward a possible SPP algorithm with duplicated genes is to understand in more details the SCJ-TD-FD directed distance algorithm.

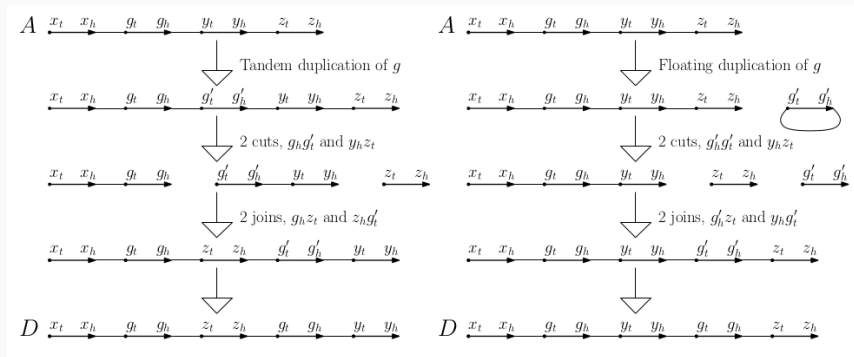
We are interested in the directed distance because in the framework of the SPP, we want to compute the distance along branches of a phylogenetic tree, so from an ancestor to a descendant.

We did see that the distance formula is very simple, actually a natural generalization of the SCJ formula with no duplications:

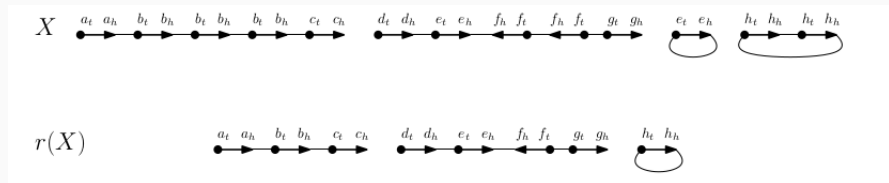
$$d_{SCJ-TD-FD}(A, D) = |A - D| + |D - A| + 2n_d$$

where n_d is the number of duplications in D .

An example of two SCJ-TD-FD scenarios



Actually, in the formula given two slides above, I did hide an assumption: I assumed that D is *reduced*, i.e. has no array of tandemly duplicated genes.



Reducing a genome X consists in iteratively removing gene copies belonging to tandem arrays (linear and circular) but if there is a unique single-gene circular chromosome.

For a genome X , we denote by $t(X)$ the number of adjacencies to remove to reduce it. For an instance (A, D) of the SCJ-TD-FD directed distance we have

$$d_{SCJ-TD-FD}(A, D) = d_{SCJ-TD-FD}(A, r(D)) + t(D)$$

If we denote by $\delta(A, r(D))$ the absolute difference in gene numbers between A and $r(D)$, this gives

$$d_{SCJ-TD-FD}(A, D) = |A - r(D)| + |r(D) - A| + 2\delta(A, r(D)) + t(D)$$

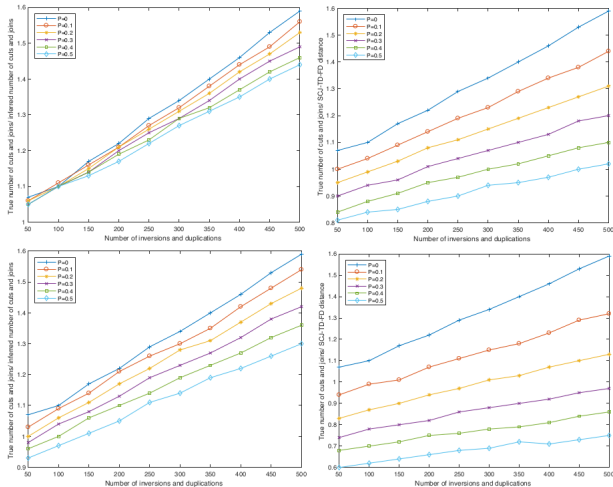
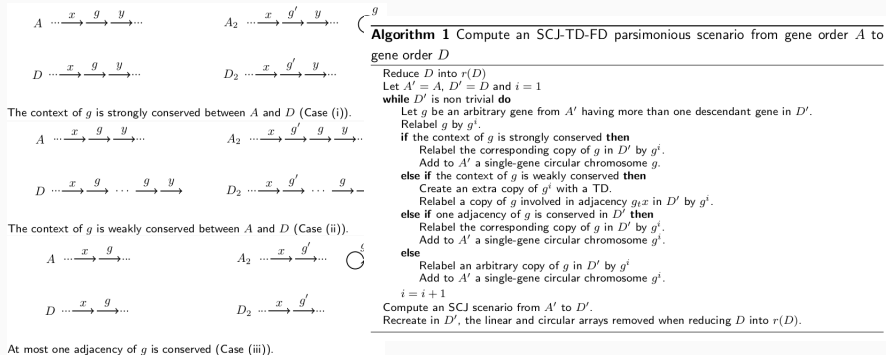


Figure 5 Experimental results, for four duplications parameters – single-gene segmental duplication (top row), two-genes segmental duplication (second row), five-genes segmental duplications (third row), variable length segmental duplications (bottom row) – and two measured quantities – inferred cuts and joins (left column) and SCJ-TD-FD distance (right column).

We can compute the distance in linear time. Can we compute a parsimonious scenario in linear time? The key point is to decide, for every duplicated copy to add, to decide if it should be added through a tandem duplication (TD) or a floating duplication (FD).



Open questions: counting the number of scenarios, sampling scenarios.

We assume we are given k reduced genomes D_i , all on the same alphabet of gene families. We want to find a trivial genome on this alphabet (i.e. a genome with a single copy of each gene) that will minimize the objective function

$$\sum_{i=1}^k d_{SCJ-TD-FD}(M, D_i)$$

We assume we are given k reduced genomes D_i , all on the same alphabet of gene families. We want to find a trivial genome on this alphabet (i.e. a genome with a single copy of each gene) that will minimize the objective function

$$\sum_{i=1}^k d_{SCJ-TD-FD}(M, D_i)$$

i.e.

$$\sum_{i=1}^k |M - D_i| + |D_i - M| + 2\delta(M, D_i)$$

We can reformulate the objective function using the classical formula for the symmetric difference

$$|M - D| + |D - M| = |M| + |D| - 2|M \cap D|$$

This gives a new formulation of the objective function to minimize

$$\sum_{i=1}^k |D_i| + 2 \sum_{i=1}^k \delta(M, D_i) - 2 \sum_{i=1}^k |M \cap D_i| + k|M|$$

This implies that actually we want to maximize

$$2 \sum_{i=1}^k |M \cap D_i| - k|M|$$

We can formulate the problem as a Maximum Weight Matching (MWM) problem.

For an adjacency a let $\gamma_i(a) = 1$ if $a \in D_i$ and 0 otherwise and let

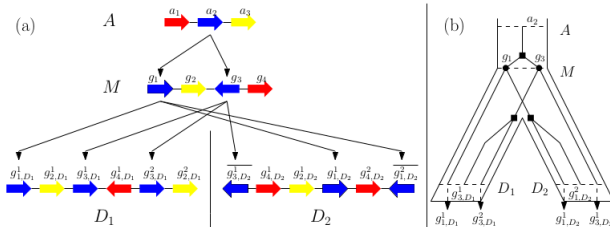
$$w(a) = 2 \sum_{i=1}^k \gamma_i(a) - k$$

Claim. Solving the MWM on the graph whose vertices are gene extremities and edges are adjacencies observed in the D_i s weighted by $w(\cdot)$ provides an optimal median M of the D_i s.

Note that actually, we do not need to consider the adjacencies with a negative weight, so we are sure that the median contains only edges seen in at least half of the D_i s.

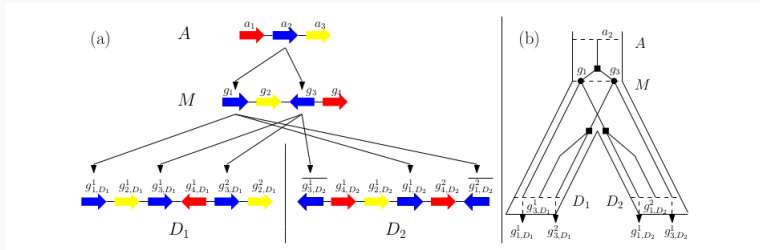
We assume we are given k reduced genomes D_i , one ancestral genome, all on the same alphabet of gene families, the gene content of a median M and orthology relations between M and the D_i s. We want to find a M that will minimize the objective function

$$d_{SCJ-TD-FD}(A, M) + \sum_{i=1}^k d_{SCJ-TD-FD}(M, D_i)$$

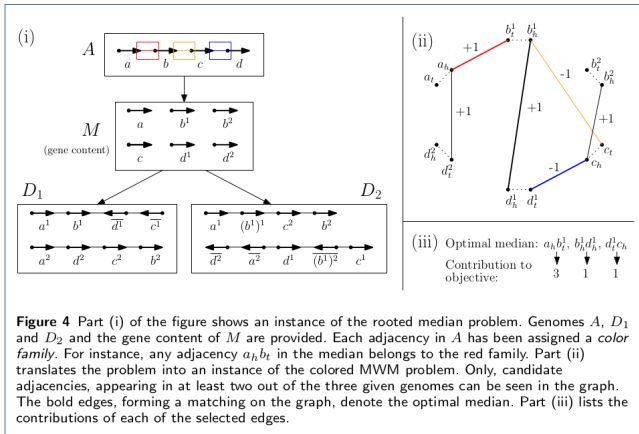


Remark 1. This is the simplest instance of an SPP problem, where we aim to find a single ancestral genome.

Remark 2. How can we assume that we already know the gene content of the median genome and orthology relations? In our context, we assume we have computed *reconciled gene trees* for our data (see block 1 seminar on gene families).



The rooted SCJ-TD-FD median is NP-hard (reduction to a variant of 3-SAT), but can be formulated as a colored matching problem and solved efficiently with an ILP.



Events	Adj. in true median	Cand. adj.	Adj. in ILP median	Precision	Recall	% Adj. common to all genomes	% Adj. common to two genomes	No. of optimal solutions	Avg. time per run (in sec)
100	1514	1503	1493	0.9998	0.9859	86.43	13.57	2.3	53
200	1107	1062	1044	0.9991	0.9428	69.49	30.51	15.8	29
300	1312	1192	1155	0.9985	0.8758	52.94	47.06	40.3	38
400	1151	985	961	0.9981	0.8329	49.44	50.56	393.7	51
500	1430	1174	1132	0.9972	0.7897	46.68	53.32	3682.6	84

Table 1 Statistics of the ILP median experiment on simulated data.

We did see that the rooted median is intractable, so we know that the SCJ-TD-FD SPP is also intractable.

Nevertheless, there is a relatively simple ILP for the rooted median, so we can ask if we can extend it to the SPP.

Let's first formulate the SPP.

We have a phylogenetic tree, with a gene order at each leaf.

For each ancestral node of the tree, we have the gene content of the corresponding ancestral species. E.g. from reconciliations.

Moreover, for each branch of the tree, we have the orthology relations along this branch. Again, this can originate from reconciled gene trees.

Last, we assume that for each ancestral node, we have a list of *candidate adjacencies*, each with a weight w between 0 and 1. They come from the DeCoStar algorithm we will describe later.

Given all the input described in the previous slide, we aim to select, at each ancestral node, a subset of candidate adjacencies that form a valid gene order and minimizes the following objective function

$$\sum_{v \in V} \left(\alpha \left(\sum_{a \in v} (1 - p_{v,a}) w_{v,a} \right) + (1 - \alpha) d_{DSCJ}(u, v) \right)$$

where d_{DSCJ} is for $d_{SCJ-TD-FD}$, $w_{v,a}$ is the weight of candidate adjacency a at node v and $p_{v,a}$ is 1 if adjacency a is selected at node v and 0 otherwise.

The SCJ-TD-FD SPP can be solved, for mixed ancestral gene orders, by a relatively simple minimization ILP.

$$\sum_{v \in V} \left\{ K_v + (1 - \gamma) \left(\sum_{a_{par} \in u} c_{u,v,a_{par}} + \sum_{a_{par} \in u} c_{v,u,a_{par}} + 2\delta(u,v) - 2t(v) + 2 \sum_{g \in G_u} \alpha_g \right) \right\}$$

subject to:

$$c_{u,v,a_{par}} \geq 0 \quad \forall a_{par} \in u \quad (1)$$

$$c_{u,v,a_{par}} \geq p_{u,a_{par}} - \sum_{a \in F_a} p_{v,a} \quad \forall a_{par} \in u \quad (2)$$

$$c_{v,u,a_{par}} \geq 0 \quad \forall a_{par} \in u \quad (3)$$

$$c_{v,u,a_{par}} \geq \sum_{a \in F_a} p_{v,a} - p_{u,a_{par}} \quad \forall a_{par} \in u \quad (4)$$

$$\sum_{a=(x_t,y)} p_{v,a} \leq 1 \quad \forall x, \forall v \quad (5)$$

$$\sum_{a=(x_h,y)} p_{v,a} \leq 1 \quad \forall x, \forall v \quad (6)$$

The constraints (6.1-6.4) ensure that the values for $c_{u,v,a_{par}}$ and $c_{v,u,a_{par}}$ are chosen correctly according to the $p_{v,a}$ values. The constraints (6.5-6.6) ensure the consistency of the genome. In other words, they ensure that each extremity takes part in at most one adjacency.

Comment 1. The result of the ILP is a set of mixed gene orders. If we want to have linear gene orders, we need to add extra constraints forbidding explicitly the obtained circular fragments and repeat. In practice, this works well.

Comment 2. In practice, the objective function is highly dominated by cuts and joins needed to handle gene losses along the tree. Experimental results are good when there is no gene loss and become quickly very bad when gene losses are allowed.

We have now a full pipeline to handle the SPP in a context of duplicated genes, using the SCJ-TD-FD model and DeCoStar.

Unfortunately, it does not work great because of the combination of the model being not robust to gene losses and the mosquito data we used it on containing a lot of artefactual gene losses due to errors in clustering genes into families.