

THE PROBLEMS OF RECONCILING GENE AND SPECIES TREES, MAPPING A GENE TREE INTO A SPECIES TREE, AND GENE TREE INFERENCE

K.Yu. Gorbunov*, V.A. Lyubetsky

Institute for Information Transmission Problems of the Russian Academy of Sciences (Kharkevich Institute), 19 Bolshoy Karetny, 127994 Moscow, Russia

Keywords: phylogenetics, fast algorithms, tree inference, trees reconciliation, supertree, gene tree mapping, cost of mapping, gene tree inference.

Track: Open Problems

*Corresponding author: gorbunov@iitp.ru

1. Reconciliation of a set of gene trees. A long recognized problem is inference of a tree S that reconciles a set of input trees G_i , with leaves in each G_i being assigned homologous sequences from an i -family (homologous genes or regulatory regions with or without the regulated gene, etc.). Usually the tree S is a tree of species or other taxonomic units. Assume leaves in S are labeled with species names, leaves in G_i – with pairs of gene-species names; paralogs are allowed. We further develop a traditional approach to find the tree S such that it minimizes the total cost of mappings of individual trees into S , [1]. Let us call S a supertree, each G_i – a gene tree, each sequence – a gene. A mapping of G_i into S implies a fixed set of evolutionary events. The *standard set* contains only gene duplications and losses. The *extended set* additionally contains horizontal gene transfers, gains, etc. The *total cost* is the sum of “individual” costs for all G_i mappings into the tree S , and is similar to the cost from [1] in case of the standard set. Is computing total/individual costs possible in small polynomial time? This question is tackled below.

Under a traditional approach, the supertree building problem is NP-hard, i.e., any algorithm to solve it correctly must possess exponential complexity. Numerous heuristics exist, but they generally do not find the global minimum of the total mapping cost. We proposed a reformulation of this problem that allows a computationally effective deterministic algorithm and meets many biological prerequisites. Namely, the supertree S is sought for such that it contains the majority of clades from input trees G_i . With the standard event set, the algorithm is mathematically correct and possesses the running time of $O(n^3 \times m^3)$, where n is the number of gene trees, and m is the total number of species [2]. For simplicity, here we assume that the average number of leaves in input trees is a multiple of m . With an extended event set, the algorithm is heuristically correct and cubic in complexity. The authors are unaware of analogous approaches in published literature. A relevant biological discussion of our approach is provided in [3]. Problem 1: is a correct inference of the supertree possible in similar or polynomial time with the extended event set?

2. Reconciliation of gene and species trees. Edges in S may be broken by inserting additional nodes, thus formally producing another tree S_0 , with nodes producing two descendant edges or only one edge. It imposes time slices such that horizontal transfers are allowed only within one slice, see [4]; in particular cases $S=S_0$. With the extended event set, we developed an algorithm that reconciles any gene tree G and S_0 , i.e., correctly computes the mapping of G into S_0 and its cost in time $O(|G| \times |S_0|)$, which gives $O(|S|^3)$. Here $||$ is the number of nodes in a tree. A mathematical proof is given in [4, 5] (refer also to a later study [6], which has used [4, 5]). Refresh that a phylogenetic net of genes or species is an acyclic directed graph with one vertex (the “root”) that can be connected by a path with any other node, and terminal vertices (the “leaves”); the leaves are labeled with species names or species-gene name pairs. An important special case is a binary net, where for each node, except for the root, one of the following is true: the node possesses only one incoming edge and no outgoing edges, or two outgoing and one incoming edge, or two incoming and one outgoing edge. In a species net, introducing time slices ordinarily generates edges with one incoming and one outgoing edge. The definitions and costs of

mapping of a gene net into a species net and mapping of a gene tree into a species tree are identical. Problem 2: for phylogenetic nets, is a correct computation of the mapping and its cost possible in the same or other polynomial time with the standard or extended set?

3. Mapping a gene tree to a species tree and its cost. Thus, Problems 1-2 are reduced to the extended event set case. We do not know their solutions in general case, but the principal question remains: what is a mapping of gene tree G into the tree S and its cost in the case of extended event set (Problem 3)? We proposed a possible definition in [4, 5] and formulate the idea below; informal motivation can be found in [4].

Below e runs over all edges of G , and d – over all edges (also referred to as *tubes*) of S_0 . The root is pictured at the top in all. A formal edge (the so-called *root edge*) enters the root from above. By definition, $e' < e$ if an edge e' is strictly below e , and $e' \leq e$ if $e' < e$ or $e' = e$. One of at most 15 evolutionary events can occur on fixed edge e in tube d ; such events are marked by index i . Two tubes d and d' belong to the same (temporal) *slice*, by definition, if the lengths of path from the root tube to d and d' are equal, correspondingly. Assume d is a terminal tube, and e is a terminal edge. If gene e belongs to species d , then the mapping $f_{\langle e, d \rangle}(e) = \langle d, \text{fin} \rangle$ and the cost $c(\langle e, d, \text{fin} \rangle) = 0$. Otherwise, if e belongs to species d' , then the mapping $f_{\langle e, d \rangle}(e) = \langle d, \text{tr}(d, d') \rangle$ and the cost $c(\langle e, d, \text{tr} \rangle) = 13$, where 13 is the cost of gene transfer without retention (all costs are conditional). Here *fin* indicates “ e belongs to d ”, *tr*(d, d') indicates that gene e is transferred from tube d to tube d' such that no copy of e remains in d . Thus, mapping $f_{\langle e, d \rangle}(e)$ is defined by induction; the basic step is defined.

Now we exemplify further induction steps $f_{\langle e, d \rangle}(e') = \langle d', \text{mark} \rangle$ where $e' < e$ and d' is in a slice later or equal then the slice d .

- 1) Assume d is a tube with single descendent tube d_1 and $\text{pass} = \arg_i \min c(e, d, i)$. Then $f_{\langle e, d \rangle}(e) = \langle d, \text{pass} \rangle$, $f_{\langle e, d \rangle}(e') = f_{\langle e, d_1 \rangle}(e')$, and $c(\langle e, d, \text{pass} \rangle) = \min_i c(\langle e, d_1, i \rangle)$. Here, the “*pass*” indicates the survival of gene e down to the next tube.
- 2) Assume d is a tube with two descendent tubes d_1 and d_2 and $\text{passl} = \arg_i \min c(e, d, i)$. Then $f_{\langle e, d \rangle}(e) = \langle d, \text{passl} \rangle$, $f_{\langle e, d \rangle}(e') = f_{\langle e, d_1 \rangle}(e')$, and $c(\langle e, d, \text{passl} \rangle) = \min_i c(\langle e, d_1, i \rangle) + 2$, where 2 is the cost of loss. Here, the “*passl*” indicates the survival of gene e into tube d_1 and loss of its copy in tube d_2 . Symmetric cases are not discussed everywhere.
- 3) Denote e_1 and e_2 two descendent edges of non-terminal edge e and $\text{passlr} = \arg_i \min c(e, d, i)$. Then $f_{\langle e, d \rangle}(e) = \langle d, \text{forklr} \rangle$, $f_{\langle e, d \rangle}(e') = f_{\langle e_1, d_1 \rangle}(e')$ if $e' \leq e_1$, and $f_{\langle e, d \rangle}(e') = f_{\langle e_2, d_2 \rangle}(e')$ if $e' \leq e_2$; $c(\langle e, d, \text{passlr} \rangle) = \min_i c(\langle e_1, d_1, i \rangle) + \min_i c(\langle e_2, d_2, i \rangle)$. Here, the “*forklr*” indicates divergence of gene e into e_1 in tube d_1 and e_2 in tube d_2 .
- 4) For non-terminal edge e and $\text{dupl} = \arg_i \min c(e, d, i)$, $f_{\langle e, d \rangle}(e) = \langle d, \text{dupl} \rangle$, $f_{\langle e, d \rangle}(e') = g_{\langle e_1, d \rangle}(e')$ if $e' \leq e_1$, and $f_{\langle e, d \rangle}(e') = f_{\langle e_2, d \rangle}(e')$ if $e' \leq e_2$; $c(\langle e, d, \text{dupl} \rangle) = \min_i c(\langle e_1, d, i \rangle) + \min_i c(\langle e_2, d, i \rangle) + 3$, where 3 is the cost of duplication. Here, the “*dupl*” indicates the duplication of gene e within tube d into e_1 and e_2 .
- 5) For non-terminal edge e and $\exists d' [\text{trl}(d, d') = \arg_i \min c(e, d, i)]$, where d' be a tube in the same slice with d that differs from d and minimizes the value $\min_i c(\langle e_1, d', i \rangle)$ over d' . Then $f_{\langle e, d \rangle}(e) = \langle d, \text{trl}(d, d') \rangle$, $f_{\langle e, d \rangle}(e') = f_{\langle e_1, d' \rangle}(e')$ if $e' \leq e_1$, and $f_{\langle e, d \rangle}(e') = f_{\langle e_2, d \rangle}(e')$ if $e' \leq e_2$; $c(\langle e, d, \text{trl}(d, d') \rangle) = \min_i c(\langle e_1, d', i \rangle) + \min_i c(\langle e_2, d, i \rangle) + 11$, where 11 is the cost of transfer with retention. Here, the “*trl*(d, d')” indicates that gene copy e_1 is transferred from tube d into tube d' , and gene copy e_2 remains in d . *Etc.*

4. Gene tree reconstruction. Problem 4 of inferring a gene tree from a multiple protein alignment has no rigorous solution, like the described above, even in important special cases, and thus always relies on heuristics. However, it is known to be NP-hard and can be formulated in terms of maximizing a defined functional. Therefore, a correct algorithm of polynomial (especially low) complexity to solve the problem is possible only after its reformulation. Importantly, such a reformulation was biologically relevant.

The authors propose the problem restatement and a polynomial algorithm that correctly infers the gene tree in special cases. In simulations of a general case the algorithm was shown to be very fast and in about 75% cases reconstruct a tree very close to that produced by PhyloBayes v.3.3 in much longer time.

In the restatement, the tree is sought for among trees consisting of clades from a prebuilt set P that possesses the following property: each set in P can be split into two subsets also from P , and so on until singlet sets are obtained that correspond to alignment rows (ref. to further as rows). The solution found by the algorithm depends heavily on the “correctness” of set P .

In the algorithm, selected “non-informative” columns of the initial alignment are omitted as in [7], dynamic programming is then applied to the refined alignment for a maximum likelihood inference of one tree that contains all clades from P per alignment column. The obtained set of trees is reconciled using a supertree building algorithm described in section 1 and in [2, 3].

5. Construction of set P . Define the row length as the number of amino acid residues. Define the length of a rows pair as the number of non “gap-to-gap” columns. Let X be any set of rows that differs from set X of all alignment rows. The set cardinality is designated $|X|$. Define $I(X)$ the sum of self-similarities of each row from X divided by the total length of all rows from X . Define $S(X)$ the sum of similarities of all row pairs from X divided by the total pairs length and $I(X)$; analogously define $S(X)$. Define $m(X)$ the minimal value of S over all pairs from X ; $M(X)$ its maximal value over all pairs that contain one row from X and another – from its complement. Let $l(m)$ be a linear function, where m decreases from $|X|$ to 2, $l(|X|) = S(X)$, and the derivative equals α , the algorithm parameter.

Set P is constructed with induction. All singlet sets (i.e. consisting of a single row) are included in P ; at each induction step, if sets X_1 and X_2 from P do not intersect, and for their union X the conditions $S(X) \geq l(|X|)$, $m^2(X) + \beta > M^2(X)$ are true then X is included in P . Parameters α and β are fitted, and in this study typically were $\alpha = 0.002556$, $\beta = 0.390625$. Although this algorithm is also heuristic, Problem 5 implies a proof of the above mentioned statements, which might be expected to achieve because the proof was obtained for Problem 1 that differs from Problem 5 by a single transparent condition.

References

- [1] R. Guigo, I. Muchnik, T.F. Smith. Reconstruction of ancient molecular phylogeny. *Molecular Phylogenetics and Evolution* **6** (1996), 189–213.
- [2] K.Yu. Gorbunov and V.A. Lyubetsky. The tree nearest on average to a given set of trees. *Problems of Information Transmission* **47**(3) (2011), 274–288.
- [3] K.Yu. Gorbunov and V.A. Lyubetsky. Fast Algorithm to Reconstruct a Species Supertree from a Set of Protein Trees. *Molecular Biology (Mosk)*, **46**(1) (2012), 161–167.
- [4] K.Yu. Gorbunov and V.A. Lyubetsky. Reconstructing the evolution of genes along the species tree. *Molecular Biology (Mosk)* **43**(5) (2009), 881–893.
- [5] K.Yu. Gorbunov and V.A. Lyubetsky. An algorithm of reconciliation of gene and species trees and inferring gene duplications, losses and horizontal transfers. *Information Processes* **10**(2) (2010), 140–144 (in Russian).
- [6] J.-P. Doyon, C. Scornavacca, K.Yu. Gorbunov, G.J. Szeollosi, V. Ranwez, V. Berry. An Efficient Algorithm for Gene/Species Trees Parsimonious Reconciliation with Losses, Duplications and Transfers. *Lecture Notes in Bioinformatics*, **6398** (2010), 93–108.
- [7] V.A. Lyubetsky, K.Yu. Gorbunov, V.V. Vyugin, L.Yu. Rusin. Removing noise in a multiple protein alignment. *Information processes*, **5**(5) (2005), 380–391 (in Russian).