

COMENIUS UNIVERSITY IN BRATISLAVA  
FACULTY OF MATHEMATICS, PHYSICS AND INFORMATICS

**SOFTWARE FOR ISOMETRIC GENE  
TREE RECONCILIATION**

Master's thesis

COMENIUS UNIVERSITY IN BRATISLAVA  
FACULTY OF MATHEMATICS, PHYSICS AND INFORMATICS

# SOFTWARE FOR ISOMETRIC GENE TREE RECONCILIATION

Master's thesis

Study programme: Applied Computer Science  
Field of study: Computer Science  
Department: Department of Computer Science  
Supervisor: doc. Mgr. Bronislava Brejová, PhD.



Univerzita Komenského v Bratislave  
Fakulta matematiky, fyziky a informatiky

---

## ZADANIE ZÁVEREČNEJ PRÁCE

**Meno a priezvisko študenta:** Dominika Mihálová  
**Študijný program:** aplikovaná informatika (Jednoodborové štúdium, magisterský II. st., denná forma)  
**Študijný odbor:** informatika  
**Typ záverečnej práce:** diplomová  
**Jazyk záverečnej práce:** anglický  
**Sekundárny jazyk:** slovenský

**Názov:** Software for isometric gene tree reconciliation  
*Softvér pre izometrickú rekonciliáciu génových stromov*

**Anotácia:** Izometrická rekonciliácia génových stromov je výpočtový problém, v ktorom je cieľom identifikovať zodpovedajúce si vrcholy v dvoch stromoch, z ktorých jeden reprezentuje evolučnú históriu skupiny organizmov a druhý reprezentuje evolučnú históriu jedného génu v rámci týchto organizmov. Cieľom práce je implementovať praktický softvér na rekonciliáciu génových stromov a experimentálne otestovať jeho presnosť na simulovaných aj reálnych biologických dátach.

**Vedúci:** doc. Mgr. Bronislava Brejová, PhD.  
**Katedra:** FMFI.KI - Katedra informatiky  
**Vedúci katedry:** prof. RNDr. Martin Škoviera, PhD.  
**Dátum zadania:** 14.10.2019

**Dátum schválenia:** 29.11.2019  
prof. RNDr. Roman Ďurikovič, PhD.  
garant študijného programu

---

študent

---

vedúci práce

# Abstrakt

Cieľom diplomovej práce bolo implementovať softvér pre izometrickú rekongiliáciu génových stromov a experimentálne otestovať jeho presnosť na simulovaných aj reálnych biologických dátach. Práca je rozdelená do piatich kapitol.

Prvá kapitola je venovaná základnej terminológii z oblasti bioinformatiky a prehľadu rôznych prístupov k riešeniu problému rekongiliácie s názornými riešeniami tejto problematiky v podobe softvérov.

Ďalšia časť uvádza algoritmy a implementovaný softvér. Opisujú sa vytvorené algoritmy, špecifikujú sa vlastnosti implementovaného softvéru, spôsob spracovania vstupov a následné výstupy.

Štvrtá kapitola prezentuje testovaciu sadu a opisuje vykonané experimenty na implementovanom softvéri pre izometrickú rekongiliáciu.

Záverečná kapitola sa zaoberá interpretáciou výsledkov testovania implementovaného softvéru.

**Kľúčové slová:** izometrická rekongiliácia génového stromu, nepresné dĺžky hrán, fylogenetický strom

# Abstract

The main goal of the diploma thesis was to implement software for isometric gene tree reconciliation and to experimentally evaluate its accuracy on simulated and real biological data. The thesis is divided into five chapters.

The first chapter presents the basic terminology in the field of bioinformatics and an overview of different approaches to solving the problem of reconciliation with concrete solutions to this problem in the form of software.

The next section presents algorithms and implemented software. The created algorithms are described, the properties of the implemented software, the method of input processing and subsequent outputs are specified.

The fourth chapter presents a test set and describes the experiments performed on the implemented software for isometric reconciliation.

The final chapter deals with the interpretation of the results of testing the implemented software.

**Keywords:** isometric gene tree reconciliation, inexact branch lengths, phylogenetic tree

# List of Figures

1.1	Unrooted tree and its rooted version . . . . .	4
1.2	Reconciliation and evolutionary history . . . . .	6
1.3	Isometric reconciliation . . . . .	12
1.4	Isometric reconciliation with inexact branch lengths . . . . .	13
2.1	Rooting the gene tree: special case condition . . . . .	20
3.1	Entity-relationship diagram with differences . . . . .	28
4.1	Duplication consistency score . . . . .	33

# Contents

<b>Introduction</b>	<b>1</b>
<b>1 Overview</b>	<b>2</b>
1.1 Background	2
1.2 Different approaches to gene tree reconciliation	4
1.2.1 Scoring gene tree reconciliation	6
1.2.2 Probabilistic gene tree reconciliation	9
1.2.3 Isometric gene tree reconciliation	11
1.2.3.1 Exact branch length	11
1.2.3.2 Inexact branch lengths	12
<b>2 Algorithms</b>	<b>17</b>
2.1 Rooting the gene tree	17
2.2 Counting algorithm	19
2.2.1 Preprocessing	21
2.2.2 Main algorithm	24
<b>3 Implementation</b>	<b>26</b>
3.1 Classes and variables	26
3.2 Differences from the original source code	27
<b>4 Experiments</b>	<b>32</b>
4.1 Simulated dataset	32
4.2 Real dataset	33
<b>5 Results</b>	<b>34</b>
<b>Conclusion</b>	<b>35</b>

<b>Appendix</b> . . . . .	<b>36</b>
---------------------------	-----------



# Introduction

# 1 Overview

In this chapter, we introduce basic information and define essential terminology from the field of bioinformatics. We describe different ways for solving the problem of gene tree reconciliation with examples of existing software in more details.

## 1.1 Background

Every organism has its complete set of genetic information encoded in a genome. A genome consists of several DNA (deoxyribonucleic acid) molecules and contains all the information, which are required for the organism to function.

DNA is a long molecule composed of two complementary strands. Each strand is made up of four chemical bases: adenine, guanine, cytosine and thymine, and connected to its complementary strand by pairing rules, where adenine is paired up with thymine and cytosine is paired up with guanine. The sequence of these bases encodes the genetic information important for building and maintaining an organism. Specific parts of DNA are called genes.

A gene is a subsequence of a DNA strand that contains information for the synthesis of a specific molecule, usually a protein. It is a basic unit of heredity.

The DNA sequence of a gene can be altered by mutation. It is a process that allows small changes in the DNA of organisms that refers to differences between individuals within a population. We will be working with two types of mutations: duplication and gene loss.

Duplication is a type of mutation where one or more genes are copied and inserted to some other position in the same genome. A duplicated gene sometimes develops a new function [9].

The opposite of duplication is gene loss (deletion). It is a type of mutation in which

some part of a DNA sequence containing a gene is left out from the genome during DNA replication or it loses its function.

Speciation is an evolutionary process of in which a single population evolves populations into two distinct species. It can happen for various reasons, for example, when a group separates from other members of its species to a different geographical area. Members of a new group develop their own unique characteristics due to the demands of another environment and this process will differentiate the new species.

Duplications and speciations result in the formation of groups of similar genes, called gene families, from a single gene. A gene family consists of evolutionarily related genes from one or multiple species, which are structurally and usually functionally similar.

Evolutionary relationships formed by evolutionary events are represented in a form of graph called a phylogenetic tree, which is a branching diagram that shows evolutionary relationships between organisms. A phylogenetic tree is a tree  $T$  with nodes  $V(T)$ , edges  $E(T)$  and leaves  $L(T)$ . It is called weighted when branch length  $w(u, v)$  is defined for each edge  $(u, v)$ .

Phylogenetic trees can be either rooted or unrooted (Fig. 1.1). A rooted tree is a phylogenetic tree  $T$  where for  $(u, v) \in E(T)$ : node  $u$  is the parent of node  $v$ , node  $v$  is the children of node  $u$ ,  $root(T)$  does not have parent and leaves  $L(T)$  do not have children. An ancestor of node  $v$  is any node of tree  $T$  on the path from node  $v$  to  $root(T)$ . Every ancestor have at least one descendant. A descendant of node  $v$  is any node of tree  $T$  of which  $v$  is an ancestor [12]. We will denote for  $u, v \in V(T)$  that  $v <_T u$  if  $u$  is ancestor of  $v$  and  $v$  is descendant of  $u$ .

Every rooted tree has a height, which symbolizes the longest path. It is the number of nodes between the  $root(T)$  and one of the leaves.

Nodes in rooted trees have levels. The level of node  $l(u)$  is the number of ancestors on the path from the node  $u$  to the  $root(T)$ . The root of a tree  $root(T)$  has level 0, since it has no parents.

If a rooted tree is weighted, nodes in the tree have depths. The depth of node  $u$ ,  $D(u)$ , is the sum of the lengths of all edges between node  $u$  and the  $root(T)$  of a rooted tree.

For a group of nodes in a rooted tree, their lowest common ancestor (LCA) is

the farthest node from the root that has all nodes in the group as descendants.

An unrooted tree is a phylogenetic tree without root. Unrooted tree can be rooted by placing a root  $r$  on some edge  $(u, v)$ . The original edge  $(u, v)$  is subdivided into two edges  $(u, r)$  and  $(r, v)$ . If edge  $(u, v)$  is weighted, then  $w(u, r) + w(r, v) = w(u, v)$ .

A special type of an unrooted tree is a semi-rooted tree, where we presume the new root of an unrooted gene tree  $G$  is positioned on edge  $(u, v) \in E(G)$  and subtrees are rooted at nodes  $u$  and  $v$ .

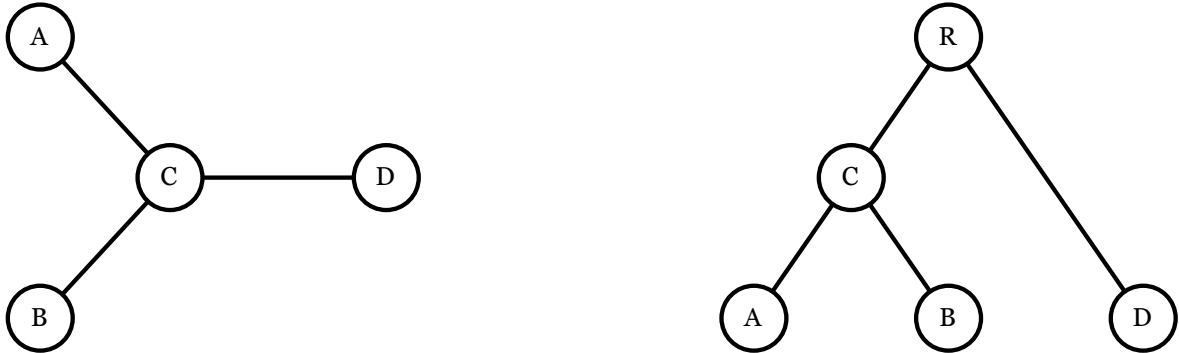


Figure 1.1: Unrooted tree(left) and its rooted version (right) tree. We placed the root  $R$  on the edge  $(C, D)$  replacing it with two new edges  $(C, R)$  and  $(R, D)$ .

We will be using two types of phylogenetic trees for showing evolutionary relationships: species trees (to describe the evolution of a set of species) and gene trees (to describe the evolution of a particular gene).

A species tree is a phylogenetic tree  $S$  where  $L(S)$  represent present-day species and internal nodes from  $V(S)$  represent speciation events in the history.

A gene tree is a phylogenetic tree  $G$  where  $L(G)$  represent present-day copies of the gene and internal nodes from  $V(S)$  represent duplication and gene loss events in the history.

Phylogenetic trees are reconstructed from a multiple alignments of the DNA sequences of present-day species by various methods [10] to find the most likely phylogenetic tree to given DNA sequences.

## 1.2 Different approaches to gene tree reconciliation

Evolutionary history is a possible sequence of evolutionary events that lead to observed members of a gene family in present-day species (Fig. 1.2). It illustrates how

many duplications and gene losses happened during the evolution of one or more genes inside the evolution of a group of species.

The problem of gene tree and species tree reconciliation was introduced in 1979 by Goodman et al. [11] as a method to infer the evolutionary history of duplications and gene losses in a gene family to decode evolutionary relationships between copies of a gene. The goal of reconciliation consists in mapping nodes of a gene tree into a species tree and thus inducing the evolution of a gene family in terms of speciations, duplications and gene losses. An important prerequisite for reconciliation is to have a gene tree without errors as misplaced leaves can lead to a different history of the gene family.

**Definition 1** *A reconciliation between gene tree  $G$  and species tree  $S$  is mapping  $\phi : V(G) \rightarrow V(S)$  such that:*

1.  $\forall u \in L(G) : \phi(u) = \mu(u)$
2.  $\forall u, v \in V(G)$  such that  $v <_G u$ :  $\phi(v) <_S \phi(u)$

An example of gene tree reconciliation is shown in Figure 1.2. We are given the gene tree  $G$ , the species tree  $S$  and a leaf mapping  $\mu : L(G) \rightarrow L(S)$  that maps each leaf from  $G$  to leaf of its species in  $S$ . We will map internal nodes according to the second condition in Definition 1. The node  $d$  is mapped to node  $Y$ , because it has  $\phi(d)$  and  $\phi(c)$  as descendants. It cannot be mapped to node  $X$  since node  $X$  does not have  $\phi(c)$  as descendant thus  $\phi(c) <_S \phi(d)$  would not hold. Then node  $e$  is mapped above node  $Y$  to have  $\phi(a)$  and  $\phi(d)$  as descendants.

The LCA-mapping  $\sigma : V(G) \rightarrow V(S)$  maps each node  $u \in V(G)$  as low as possible to the unique node  $\sigma(u) = LCA(\mu(v) | \forall v \in L(G), v <_G u)$  in  $S$ . It satisfy both conditions in Definition 1 and minimize the number of duplications and gene losses. This reconciliation can be found in linear time [12].

In this work, we divide approaches to reconciliation into three types: scoring, probabilistic and isometric. Every one of these approaches has its way to compute the gene tree reconciliation. They will be described with presented examples of software that have implemented gene tree reconciliation.

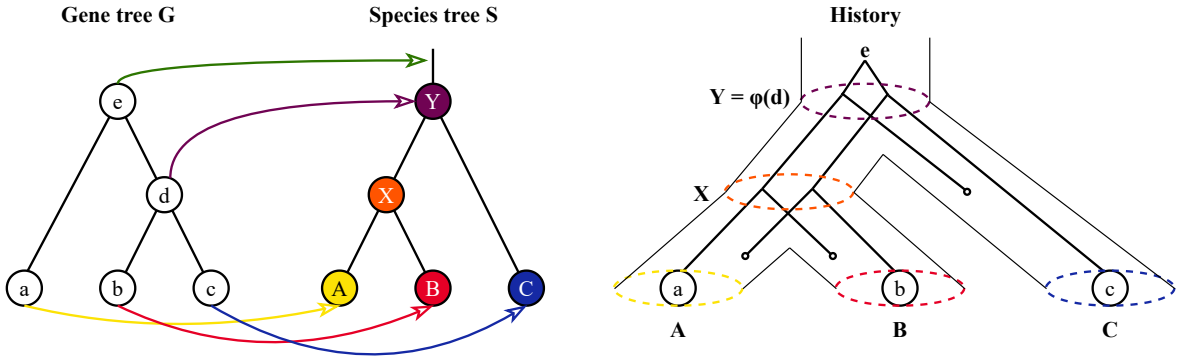


Figure 1.2: Reconciliation and evolutionary history. On the left, the gene tree  $G$  is mapped to the species tree  $S$ . On the right, we can see the evolutionary history implied by this reconciliation. This history contains one duplication  $e$ , two speciations  $Y$  and  $X$  and three gene losses (empty circles).

### 1.2.1 Scoring gene tree reconciliation

One of the known approaches to find the best reconciliation is to minimize the duplication-loss score, which signifies the sum of duplications and gene losses during the reconciliation. Various software for scoring gene tree reconciliation are known such as TreeBeST, TreeFix, Treerecs or Notung.

#### TreeBeST

Software TreeBeST [13] takes a rooted species tree and multiple sequence alignment of gene trees for the gene family as input. It uses a method to merge various input gene trees into one gene tree with a model to penalize duplications and gene losses relative to a known species tree. The method first resolves the topology of a gene tree with five methods: neighbour-joining synonymous distance, non-synonymous distance, p-distance and max-likelihood under the WAG and the HKY model. Then the topology is bootstrapped 100 times. The output of the software is one rooted tree.

The TreeBeST reconciliation method was then compared with PhyML+RAP method, where multiple sequence alignment of gene trees are reconciled to one gene tree without the presence of species tree [19]. To evaluate these two methods, authors developed a duplication consistency score represented by  $\frac{\text{intersection}}{\text{union}}$  of species between left and right branches. Low duplication consistency score means poor topology of the resultant gene tree. PhyML+RAP approach leads to more duplication nodes than TreeBeST and the supplementary duplication nodes made by PhyML+RAP had a low

duplication consistency score. Authors found this result unexpected as TreeBeST uses the species tree and tends to produce duplication when the gene tree has extensive extant members on each side of the duplication.

### **TreeFix**

TreeFix [1] takes a rooted species tree, a maximum likelihood gene tree and a multiple sequence alignment of gene trees for the gene family as input. It infers duplications and gene losses using maximum parsimony reconciliation with a duplication-loss cost function, which looks for the reconciliation with the minimum total number of duplications and gene losses. Duplication ( $D$ ) and gene loss ( $L$ ) have their costs ( $c_D$  and  $c_L$ ) that are set to one by default and can be changed by the user. The duplication-loss cost for one reconciliation can be written as:  $c_D \cdot D + c_L \cdot L$ .

The main method [20] uses a hill-climbing search strategy to find an optimal rooted gene tree with the statistically equivalent likelihood to the given maximum likelihood gene tree and with minimum duplication-loss cost as output. Basics of the method are to compute duplication-loss cost and perform neighbour interchange and subtree prune and regraft on the current optimal gene tree. After finding proposals, it chooses the proposal with the lowest duplication-loss cost and with the statistically equivalent likelihood to the given maximum likelihood gene tree. This process is repeated for a given number of iterations.

Authors compared TreeFix with RAxML, SPIMAP, TreeBeST, Notung and tt using simulated and real datasets of two types of species: 12 *Drosophila* and 16 fungi [21]. The software was judged from the point of view of phylogenetic accuracy in 5 categories (topology, branch, orthologs, duplications, losses) and runtime. TreeFix and SPIMAP show the best accuracy in all 5 categories, Notung has slightly worse accuracy in reconstructing the topology of fungi and precision of inferring duplication and gene losses. tt has problems in the same categories as Notung. The worst accuracy demonstrate RAxML and TreeBeST. While TreeFix and SPIMAP have great phylogenetic accuracy their average running time is longer than others. The best runtime has Notung followed by tt and RAxML.

### **Treerecs**

Software Treerecs [14] takes a rooted species tree, one or more rooted or unrooted gene trees and a mapping of gene tree leaves to species tree leaves. It provides recon-

ciling gene tree within the associated species tree minimizing the duplication and gene loss score and rooting the gene tree along the way if needed. The output is, depending on the input, one or more rooted gene trees.

The main aim of the authors was to create a more efficient software. They compared Treerecs with EcceTERA, Notung and Ranger-DTL in 3 categories: root (find the root that minimizes duplication-loss score), correction and root+correction (do both previous categories at the same time). In the first category, Treerecs is better than Ranger-DTL and Notung, which shows a large increase in execution time as the number of leaves increases. Treerecs show the best performance in the correction of trees followed by Notung, while Ranger-DTL has big execution time even with a small increase of leaves. The last category is only supported by Treerecs and EcceTERA, where Treerecs has also better performance.

### **Notung**

Notung [8] takes a rooted species tree, a rooted or unrooted gene tree and the leaf mapping as input. If the gene tree is not rooted, it can be rooted by Notung rooting mode that gives each edge a root score (weighted sum of duplications and gene losses). Apart from a reconciliation of binary trees, Notung can reconcile binary gene trees with non-binary species trees and non-binary gene trees with the binary species tree. The non-binary tree is a tree with at least one polytomy (a node with more than two children). Reconciliation of binary gene trees to non-binary species trees results into binary gene tree.

They use an algorithm, that can distinguish between duplication and deep coalescence (divergence, when the time of separation of two lineages precede the time of speciation) and leads to the smaller total number of duplications and losses than duplication-loss cost function used in reconciliation of two binary trees [18]. The duplication-loss cost function is the same as in TreeFix. Reconciliation of non-binary gene tree to binary species tree results into non-binary gene tree. The general approach is to convert the non-binary gene tree to binary gene tree that has minimal duplication-loss score when reconciled with the binary species tree. The resolution is then rearranged back to non-binary gene tree, where all nodes and edges not present in the original gene tree are removed and their assigned duplications and gene losses are reassigned to their polytomy.



## 1.2.2 Probabilistic gene tree reconciliation

Probabilistic methods have been designed to increase the accuracy of reconciled trees. We introduce two software tools that use probabilistic methods to reconcile gene trees: SPIMAP and Phyldog.

### SPIMAP

SPIMAP software [16] takes a rooted species tree and multiple gene sequences of species from the species tree. Normally, the gene sequences are compared and clustered according to their similarity, which results in a set of homologous gene families. Each gene family has its multiple sequence alignment that is reconstructed into gene trees and they are reconciled with the known species tree. Into this classic pipeline, SPIMAP inserts a parameter estimation model using Bayesian approach creating a new phylogenomic pipeline. It learns duplication, gene loss rates during clustering and gene, species substitution rates during the process of alignment. These parameters are then used while building and reconciling the gene tree with the known species tree. The output is a special reconciliation file format that contains gene node ID, species node ID and evolutionary event that occurred on a given node.

The parameter estimation model [17] infers duplication and gene loss rate using the birth-death process. The birth-death process is a continuous-time process that generates a gene tree according to the constant birth rate (representing duplication) and death rate (representing gene loss). After running it for a time that represents branch length, all branches that exist at the time are "surviving" and others are "extinct". If a node has no surviving descendants, it is called "doomed". Every branch has its length, which can be written as  $\frac{\text{substitutions}}{\text{site}}$  or a product of a duration of time and a substitution rate. The substitution rates signify the number of substitutions per site per unit of time. The model computes gene-specific rate (measures all rate in a tree) for every gene family and species-specific (specifies rate to given branch in the gene tree) rate for every branch.

To determine if the new phylogenomic pipeline improved accuracy, authors compare SPIMAP with PrIME-GSR, SPIDIR, MrBayes, PHYML, BIONJ, RAxML and SYNERGY. They used the same data as TreeFix: 12 *Drosophila* and 16 fungi. Firstly, they measured the average runtime. The best runtime, under 1 minute, have RAxML, MrBayes, PhyML and BionJ, which was the fastest method. With the same amount

of iterations, SPIMAP was quicker than SPIDIR or PrIME-GSR, which was the slowest method. Next, they decided to apply the duplication consistency score (used in TreeBeST) to determine method with better accuracy. The smallest number of duplication with low duplication consistency score have SPIMAP and SYNERGY. The moderate performance shows PrIME-GSR and SPIDIR. Remaining four methods have a similar number of duplication with low duplication consistency score. Lastly, they evaluate phylogenetic accuracy depending on 5 categories (used in TreeFix) for 6 above-mentioned methods (except RAxML and SYNERGY). SPIMAP has higher accuracy in every category. In the category of inferring the topology, PrIME-GSR has slightly worse accuracy while SPIDIR, MrBayes, PHYML, BIONJ shows bad accuracy in the topology of fungi dataset. The accuracy of reconstructed branches is better than the topology in every method, SPIMAP and PrIME-GSR are first two. SPIDIR, MrBayes, PHYML, BIONJ are a little worse at the sensitivity of detection orthologs in fungi dataset. They have the biggest problem with the precision of inferring the duplications in fungi dataset and losses in both datasets. The same problem has also PrIME-GSR, but only in precision of inferring losses.

### **PhylDog**

Another software is PhylDog [2] takes multiple gene alignments, a mapping between gene names and species names, and a list of species names as input. The method infers species tree, gene trees, duplication and gene loss rates by maximizing the probability of alignments overall gene families composed from the likelihood of a phylogeny given an alignment and the likelihood of the reconciliation of a gene tree with a species tree according to duplication and gene loss rates. This method uses the birth-death process and is similar to SPIMAP, but they differ in two aspects. First, while SPIMAP assumes duplication and gene loss rates to be constant for all branches in the species tree, PhylDog chooses to use a particular pair of duplication and gene loss rates to each branch of the species tree. Second, SPIMAP requires time-anchored species tree (branch length shows the amount of time between two nodes) to compute the likelihood of a gene family. Alternately, PhylDog calculates likelihood from the expected numbers of duplications and gene losses. The output is reconciled gene trees.

PhylDog was compared with TreeBeST and PhyML [3] in terms of the number

of duplications and the reconstructed ancestral genome size. PhyML has the biggest number of predicted duplication events, TreeBeST reconciled trees with a much smaller number, but still much higher than PhylDog. The same order of software was also in the number of reconstructed ancestral genome size, where both PhyML and TreeBeST have bigger ancestral genomes that lead to deeper nodes in the species tree.

### 1.2.3 Isometric gene tree reconciliation

Another variant of reconciliation is isometric gene tree reconciliation, where both species tree  $S$  and gene tree  $G$  has known branch lengths. These branch lengths are taken into account while mapping a gene tree  $G$  to a species tree  $S$ . The output of isometric reconciliation is reconciled gene tree with preserved evolutionary distances. The branch lengths of phylogenetic tree express estimated time between evolutionary events. The time can signify the actual geological time, amount of evolutionary changes that happened on the edge or expected number of substitution per site between two nodes.

#### 1.2.3.1 Exact branch length

This problem was introduced and named by Ma et al. [15] in 2008 for the first time. They defined isometric reconciliation of rooted species tree and unrooted gene trees, where all input trees have exact branch lengths. The algorithm executes all input gene trees one by one. First of all, it maps all leaves from the gene tree to leaves in the species tree. Then it takes an unmapped node, which has to be connected with at least 2 already mapped nodes, and call function, that map the unmapped node into the species tree and root the gene tree. The presented algorithm had  $O(N^2)$  running time, where  $N$  stands for the total number of nodes in the gene tree and the species tree. However, their definition of isometric reconciliation has some flaws since it does not preserve all evolutionary distances between nodes as it allows them to develop a reconciliation that does not satisfy any history of evolution.

As a result, Brejová et al. [4] later corrected and modified algorithm by Ma et al. to a more efficient algorithm with  $O(N \log N)$  running time. Their modify algorithm firstly maps every leaf from gene tree to the species tree. Next, it maps all unmapped nodes to the species tree. After all nodes are correctly mapped, the algorithm roots

the gene tree and maps the found root to the species tree. Eventually, it verifies if the reconciled tree is correct due to the definition of isometric reconciliation. They also proposed two extensions of the problem. In the first extension, they considered both input trees (gene tree and species tree) to be unrooted and designed an algorithm with  $O(N^5 \log N)$  running time. The second extension presents an algorithm, where both input trees are rooted, but gene tree branch lengths are assumed to be scaled by an unknown scaling factor.

**Definition 2** *An isometric reconciliation between gene tree  $G$  and species tree  $S$  with exact branch lengths  $w$ , which are strictly positive, is mapping  $\phi : V(G) \rightarrow V(S) \times R$  such that:*

1.  $\forall u \in L(G) : \phi(u) = [\mu(u), 0]$
2.  $\forall u, v \in V(G)$  such that  $v <_G u$ :  $\phi(v) <_S \phi(u)$  and  $w(u, v) = d(\phi(u), \phi(v))$ , where  $d$  is the length of path between  $u$  and  $v$  in reconciled gene tree.

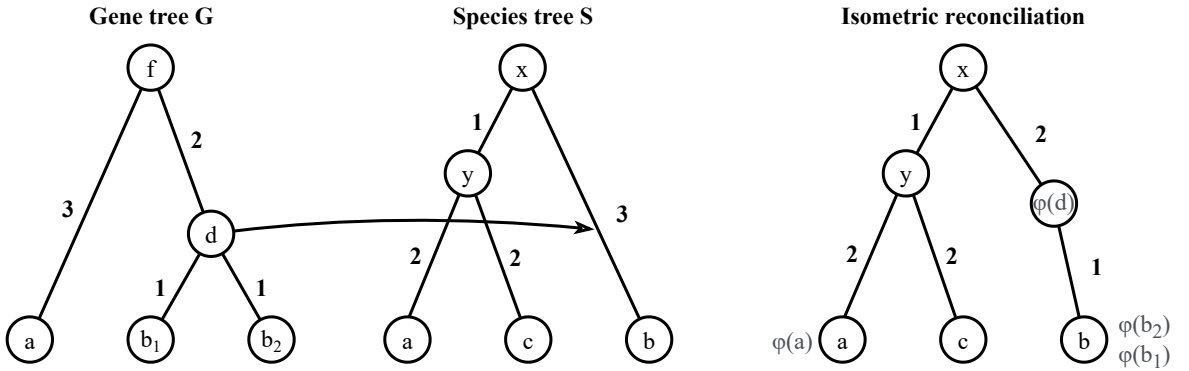


Figure 1.3: Isometric reconciliation. On the left is mapping of the node  $d$  of the gene tree  $G$  to the edge  $(x, b)$  of the species tree  $S$ . The result of the isometric reconciliation with mapped node  $d$  is on the right.

### 1.2.3.2 Inexact branch lengths

Input to the above-mentioned algorithms, gene trees and species trees with their branch lengths, are practically estimated from DNA sequences, which were gathered from present-day species [10]. The branch lengths are computed from observed mutations in collected DNA sequences. However, mutations happen randomly in the evolution and DNA sequences are also random samples from studied present-day species. It

means that inferred gene trees and species trees with their branch lengths are estimated with an error.

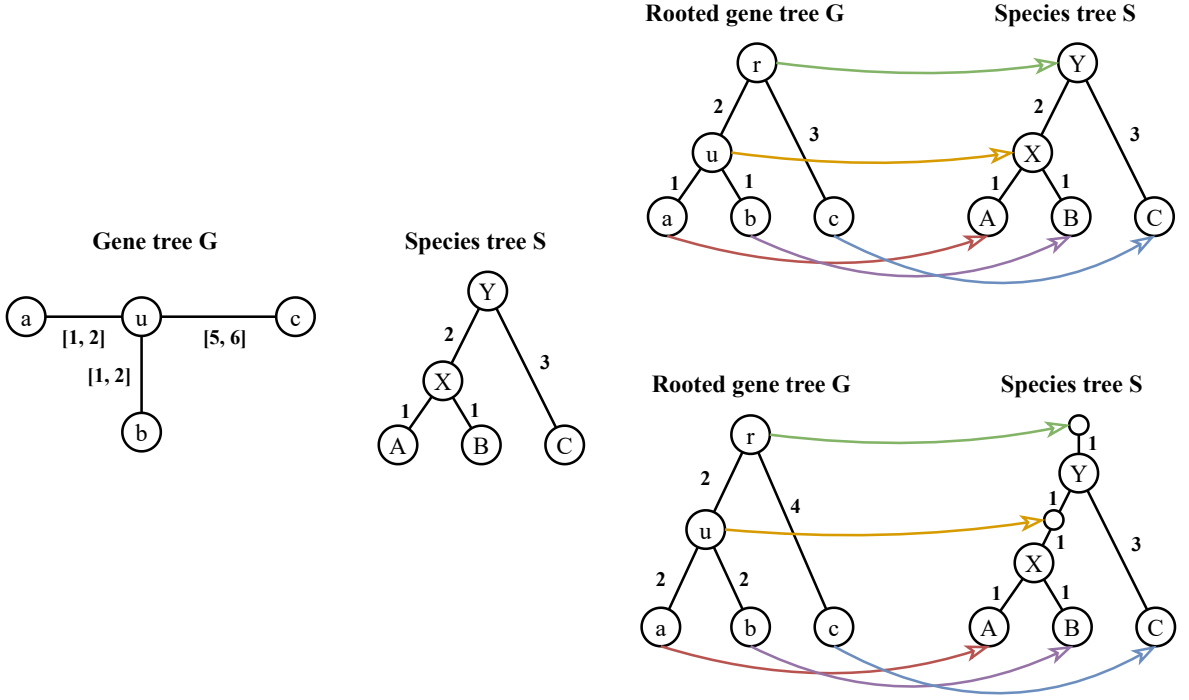


Figure 1.4: Isometric reconciliation with inexact branch lengths. On the left, we can see an unrooted gene tree  $G$  with inexact branch lengths and a rooted species tree  $S$ . On the right are two possible solutions with rooting the gene tree on the edge  $(u, c)$ : the first reconciliation is without duplications and gene losses while the second solution maps gene tree to species tree with two duplications ( $u$  and  $r$ ) and four gene losses.

To avoid defectiveness and make isometric reconciliation more precise, the isometric reconciliation with inexact branch lengths was introduced by Chládek in [6] and [5]. They define inexact branch length as an interval, where for every weighted edge stands:  $w(u, v) = \langle w(u, v)_{min}, w(u, v)_{max} \rangle$ . They present three types of algorithm.

### Linear programming algorithm

The algorithm is based on linear programming and takes a rooted species tree and a rooted gene tree with inexact branch lengths as an input. They introduce a term of mapping depth, which represents a depth of mapped node  $u \in V(G)$  in species tree:  $D(\phi(u))$ . The solution of isometric reconciliation is to find mapping depths to all nodes from gene tree to species tree. They suggest a set of 5 inequalities that can be solved by a linear program.

The first two inequalities assume that the difference of mapping depths of two nodes, where the one is a parent and the other is his child, has to be between the maximal

and minimal length of the edge. The second two inequalities are similar and say that the distance between depths of two neighbouring nodes in species tree must be within the maximal and minimal length of the edge. Last inequality restricts mapping depth of node from gene tree to be same smaller than the depth of its LCA-mapping node in the species tree. The algorithm also defines that leaves from gene tree map to leaves from species tree and the root of the species tree is in depth 0, so every node, which maps above the root has negative depths.

For  $N$  nodes, the algorithm can get to the result in polynomial time. It also works if the gene tree and species tree are non-binary trees. With a few changes in inequalities, the linear program can find reconciliation to semi-rooted and unrooted species and gene trees.

### Two-pass algorithm

The two-pass algorithm is faster than the linear programming algorithm. The input for the algorithm is a rooted species tree with exact branch lengths and a rooted gene tree with inexact branch lengths. It consists of two parts: upward and downward sweep.

The upward sweep goes from the leaves of the gene tree to the root and it computes preliminary interval  $\langle x[u]_{min}, x[u]_{max} \rangle$  for each node  $u$ . The interval is a set of all potential mapping depths values of node  $u$  over all reconciliations of a subtree of the gene tree rooted at node  $u$  to the species tree. It does not take into account the rest of the gene tree, only the descendants of node  $u$ . The highest mapping point of node  $u$ ,  $x[u]_{min}$ , is computed by taking the maximum value of both children highest mapping points subtracted by the longest branch length values. Similarly is calculated the lowest mapping point  $x[u]_{max}$ , where both children lowest mapping points are subtracted by the shortest branch length values and the minimum of these values is taken. If  $x[u]_{min} > x[u]_{max}$ , the interval is empty thus for the input is no reconciliation.

To get the final interval  $\langle X[u]_{min}, X[u]_{max} \rangle$  for each node  $u$ , the downward sweep goes from the root of the gene tree to its leaves. For all nodes in gene tree (except root), the final interval is computed from their parent final interval, where the minimal mapping depth  $X[u]_{min}$  is the maximum value of the highest mapping point of node  $u$  and minimal mapping depth of parent added by the shortest branch length value. The maximum mapping depth  $X[u]_{max}$  is computed as the minimum value of the lowest

mapping point of node  $u$  and maximal mapping depth of parent added by the longest branch length value.

Running time of this algorithm is  $O(N)$ . The same running time is for semi-rooted gene tree. In this situation, it needs to be used linear programming to obtain final intervals of the possible root and its children. The algorithm can be as well applied to an unrooted gene tree, where the running time is  $O(N^2)$ , because it needs to run on every edge as we do not know on which edge the root is located. This algorithm will be used in further work.

### **Parsimonious algorithm**

The parsimonious algorithm looks for the most parsimonious solution over all isometric reconciliations by counting the number of duplications and gene losses. The aim is to find the smallest number of duplications and gene losses. In the thesis, [5] are considered 3 different types of parsimonious algorithms for 3 types of gene tree: rooted with exact branch lengths, rooted with inexact branch lengths and semi-rooted or unrooted.

The algorithm designed for reconciliation of a rooted gene tree and species tree with exact branch lengths count duplication and gene losses on subtrees. It has the best running time of  $O(N \log N)$ . If node  $u$  from gene tree is not mapped to its LCA-mapping in the species tree, the mapping of node  $u$  to the species tree is considered as duplication. The number of gene losses is computed from a path in species tree between mappings of node  $u$  and node  $v$ , where  $u$  is the parent of  $v$ . Each node from the species tree, that occurs on this path, is considered as speciation. It creates a copy of gene represented by the edge  $(u, v)$  from gene tree, but the gene continuous to only one child, so there is a loss on the other lineage.

Improved algorithm for a rooted gene tree with inexact branch lengths and a rooted species tree with exact branch lengths counts duplication and gene losses on subintervals. It splits the mapping depth interval into non-overlapping subintervals. The goal is to compute the number of duplication and gene losses for all possible subintervals, which is done by counting function from the previous algorithm. The time complexity of the algorithm is  $O(N^3 \log N)$ .

The most parsimonious algorithm assumes semi-rooted or unrooted gene tree with inexact branch lengths and a rooted species tree with exact branch lengths. In both

cases of the gene tree, the exact location of the root is unknown. The algorithm firstly runs the previous algorithm on nodes of edge, where the possible root can be located. Then, it uses linear programming to find the location of the root based on computed numbers of duplications and gene losses in subintervals of the possible edge nodes. The running time of this algorithm is  $O(N^4 \log N)$ .



## 2 Algorithms

In this chapter, we will show different algorithms used to obtain the most parsimonious isometric gene tree reconciliation which will be implemented in our software.

The required input for our algorithms is a rooted species tree with exact branch lengths and an unrooted gene tree with inexact branch lengths. The gene tree is sequentially processed by following algorithms. Firstly, we find all possible roots and root the unrooted tree gene tree, which results in multiple rooted gene trees. Afterwards, we perform the two-pass algorithm (Chapter 1.2.3.2) and count the number of duplications and gene losses for every rooted gene tree to find the most parsimonious reconciliation.

### 2.1 Rooting the gene tree

We present an algorithm to find possible roots that are distant by given step for every edge  $e$  of an unrooted gene tree  $G$ .

For each edge  $e \in E(G)$ , we transform the unrooted gene tree into a semi-rooted gene tree by rooting the subtrees at vertices  $u \in V(G)$  and  $v \in V(G)$  of the edge  $e = (u, v)$ . The edge is subsequently used as the parameter for the Algorithm 1 to infer intervals for subdividing the given edge and completely root the unrooted gene tree  $G$ . By default, the  $\epsilon$  is set to  $1 \times 10^{-6}$  and step to 0.5.

At the beginning of the Algorithm 1, we define essential variables. Possible intervals for subdividing the edge  $e$  are stored at the variable *interval* that is also return value of the function. To cover most of the possibilities, we allow rooting the gene tree right above the vertices  $u$  and  $v$  of the edge  $e$ . Let  $r$  be the root of a semi-rooted gene tree  $G$  then we get two possible roots after subdividing the edge  $e$  right above the vertices  $u$  and  $v$ . The first possible root subdivide the edge

---

**Algorithm 1** Possible intervals to subdivide given edge  $e$ 


---

```

1: function GETINTERVALS( $e \in E(G)$ ,  $step \in R$ )
2:   totalMinLength =  $w(e)_{min}$ 
3:   totalMaxLength =  $w(e)_{max}$ 
4:   minLengthLeft = totalMinLength
5:   maxLengthLeft = totalMaxLength - step
6:   minLengthRight = 0
7:   maxLengthRight = step
8:   difference = totalMaxLength - totalMinLength
9:   differenceIsDivisibleByStep = difference / step == floor(difference / step)
10:  intervals.add( $\langle 0, \epsilon \rangle$ ,  $\langle totalMinLength, totalMaxLength - \epsilon \rangle$ )
11:  intervals.add( $\langle totalMinLength, totalMaxLength - \epsilon \rangle$ ,  $\langle 0, \epsilon \rangle$ )
12:  while  $maxLengthRight < totalMaxLength$  and  $maxLengthLeft > 0$  and
     $step \neq 0$  do
13:    intervals.add( $\langle minLengthLeft, maxLengthLeft \rangle$ ,
     $\langle minLengthRight, maxLengthRight \rangle$ )
14:    if  $difference = step$  or not  $differenceIsDivisibleByStep$  then
15:      intervals.add( $\langle minLengthRight, maxLengthRight \rangle$ ,
     $\langle minLengthLeft, maxLengthLeft \rangle$ )
16:    minLengthLeft -= step
17:    if  $minLengthLeft < 0$  then
18:      minLengthLeft = 0
19:    maxLengthLeft -= step
20:    minLengthRight += step
21:    maxLengthRight += step
  return intervals

```

---

$e$  into two edges with interval lengths  $w(u, r) = \langle 0, \epsilon \rangle$  on the left from the root and  $w(r, v) = \langle totalMinLength, totalMaxLength - \epsilon \rangle$  on the right from the root, where  $w(e)_{min}$  is original minimal length of the edge  $e$  and  $totalMaxLength$  is original maximal length of the edge  $e$ . The second option of the root subdivide the edge  $e$  with intervals of the same length, that are flipped, so the original interval  $w(u, r) = \langle totalMinLength, totalMaxLength - \epsilon \rangle$  is on the left from the root and  $w(r, v) = \langle 0, \epsilon \rangle$  in on the right from the root.

After inferring the first two options of the possible root, we run a while loop to get possible roots inside the edge  $e$ . In each iteration, we subtract the step from the minimal and maximal length of the left interval of the subdivided edge  $e$  and add a step to the minimal and maximal length of the right interval of the subdivided edge  $e$ . We added condition for special cases, where the difference between the maximal and minimal original length of edge  $e$  is equal to step or when it is not divisible by step (Fig. 2.1). With the special case condition, we can infer intervals that would be normally skipped. The while loop goes until the maximal length of the right interval is the same or bigger as the original maximal length of the edge  $e$  or the maximal length of the left interval reaches 0 or less.

Rooting the gene tree goes over all edges in the unrooted gene tree  $G$ . The running time of the function *getIntervals* in Algorithm 1 is  $O(p)$ , where  $p$  stands for number of iterations in while loop that can be expressed as  $p = \text{ceil}(totalMaxLength/step) - 1$ . The result of function *getIntervals* is a set of intervals for subdividing the edge  $e$ . Subsequently, the intervals are used for a loop to root the semi-rooted gene tree  $G$  resulting in a set of rooted gene trees with inexact branch lengths, which running time is  $O(m)$  for  $m = p + 2$ . Therefore, the total running time of rooting the unrooted gene tree  $G$  is  $O(Nm)$  with  $N$  being the number of all edges in the unrooted tree  $G$ .

## 2.2 Counting algorithm

We present an algorithm for counting the number of duplications and gene losses in a rooted gene tree  $G$  with inexact branch lengths depending on a rooted species tree  $S$  with exact branch lengths. We allow only evolutionary events as duplication, gene loss and speciation can happen in evolutionary history. The counting algorithm

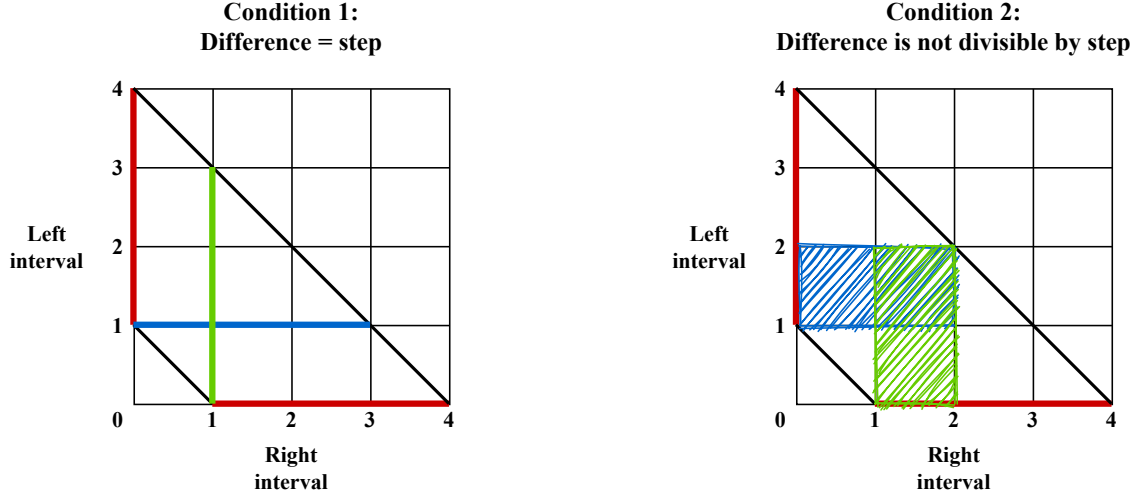


Figure 2.1: The X-axis on the graph stands for the possible length of interval on the left from the root and the y-axis signify the possible length of interval on the right from the root. Space between black lines represents all possible intervals after subdividing the edge  $e$  with the root. Red lines are for the root right above the vertices  $u$  and  $v$ , where one interval takes the original size and the second interval is between 0 and  $\epsilon$ . Blue lines mean intervals for the root inside the edge  $e$  without the special case condition. Green lines are inferred intervals for the root with special case condition

### Condition 1

Parameters:  $w(e) = \langle 1, 4 \rangle$  and  $step = 3$ .

In this case, the difference between the original minimal and maximal length is equal to the size of the step. Without the special case condition, we would only get the blue line intervals  $w(u, r) = \langle 1, 1 \rangle$  on the left from the root and  $w(r, v) = \langle 0, 3 \rangle$  on the right from root. The special case condition flip the intervals thus we get green intervals  $w(u, r) = \langle 0, 3 \rangle$  on the left from the root and  $w(r, v) = \langle 1, 1 \rangle$  on the right from the root.

### Condition 2

Parameters:  $w(e) = \langle 1, 4 \rangle$  and  $step = 2$ .

The difference between original minimal and maximal length is not divisible by the size of step, which give us only the blue lines intervals  $w(u, r) = \langle 1, 2 \rangle$  on the left from the root and  $w(r, v) = \langle 0, 2 \rangle$  on the right from root without the special case condition. After flipping the intervals with special case condition, we get green intervals  $w(u, r) = \langle 0, 2 \rangle$  on the left from the root and  $w(r, v) = \langle 1, 2 \rangle$  on the right from the root. Green and blue intervals intersect in area of  $\langle 1, 2 \rangle$  for left interval and  $\langle 1, 2 \rangle$  for right interval, but it also give us a new area of  $\langle 0, 1 \rangle$  for left interval and  $\langle 1, 2 \rangle$  for right interval.

consists of two parts: preprocessing and the main algorithm. In the preprocessing, we compute essential variables for the gene tree  $G$  that are subsequently used in the main algorithm.

The prerequisites for the counting algorithm are calculated depth  $D(a)$  and level  $l(a)$  for each  $a \in V(S)$ . For the gene tree  $G$ , we assume an interval of possible mapping

depths  $X[u] = \langle X[u]_{min}, X[u]_{max} \rangle$  for each  $u \in V(G)$  that can be computed with the two-pass algorithm (Chapter 1.2.3.2) and LCA-mapping  $\sigma(u)$  for each  $u \in V(G)$ .

### 2.2.1 Preprocessing

In the preprocessing, we calculate necessary variables for nodes in the gene tree  $G$  that are used in the main algorithm 2.2.2. For each  $u \in V(G)$  we compute  $l(u)$ ,  $speciesNodeBelow(u)$ ,  $levelDistanceFromParent(u)$  and  $mappedToLca(u)$ . The level  $l(u)$  variable of gene node  $u$  signifies number of species nodes on the path from gene node  $u$  to the root of the gene tree  $G$ . It is calculated from the  $speciesNodeBelow(u)$  variable that represents species node  $s \in V(S)$ , which is right below gene node  $u$ , so the gene node  $u$  lies on the path from species node  $s$  to  $parent(s)$ . The  $levelDistanceFromParent(u)$  means number of species nodes between gene node  $u$  and its  $parent(u)$ . It is not calculated for the root of the gene tree  $G$  as it has no parent. The variable  $mappedToLca(u)$  presents the truth value of statement:  $X[u]_{max} = D(\sigma(u))$ , thus if the node  $u$  is mapped to  $\sigma(u)$  or not.

We use three algorithms for calculating the variables: Algorithm 2 for calculating the  $speciesNodeBelow(u)$ , Algorithm 3 that sets level distance from parent  $u \in V(G)$  to its children  $\forall v \in children(u) : levelDistanceFromParent(v)$  and Algorithm 4, which uses both previous algorithms and computes the level of gene node  $l(u)$  and  $mappedToLca(u)$ .

---

**Algorithm 2** Computes species node below given gene node  $u$

---

```

1: function COMPUTESPECIESNODEBELOW( $u \in V(G), s \in V(S)$ )
2:   speciesNodeBelow = s
3:   while  $D(s) > X[u]_{max}$  do
4:     speciesNodeBelow = s
5:     if  $parent(s) = null$  then
6:       break
7:     else
8:        $s = parent(s)$ 
   return speciesNodeBelow

```

---

The *computeSpeciesNodeBelow* function in Algorithm 2 saves the given species node  $s$  to the variable *speciesNodeBelow*. It always remembers the last species node below the given gene node  $u$  and serves as a return value. The while loop iterates over species nodes on the path from given species node  $s$  to the root of the species tree. If

$D(s) > X[u]_{max}$  stands, it means the species node  $s$  is below the gene node  $u$ . Every iteration, we save the species node to the return value *speciesNodeBelow* and move on the path closer to the root by assigning  $parent(s)$  to the  $s$  variable. If the species node  $s$  does not have a parent, the gene node  $u$  is above the root of the species tree. The while loop ends, if the while condition does not stand, so the new species node  $s$  is above gene node  $u$  or if the gene node  $u$  is above the root.

---

**Algorithm 3** Sets level distance from parent to children of node  $u$

---

```

1: function LEVELDISTANCEFROMCHILDREN( $u \in V(G)$ )
2:   for  $v \in children(u)$  do
3:     levelDistanceFromParent( $v$ ) =  $l(v) - l(u) - mappedToLca(u)$ 
```

---

In the *levelDistanceFromChildren* function show in Algorithm 3, we compute the *levelDistanceFromParent( $v$ )* for each child  $v$  of the given node  $u$ . The level distance is calculated as the difference between the level of child  $v$  and the level of its parent, node  $u$ . The distance is subtracted by 1 if the node  $u$  is mapped to the same depth as its  $\sigma(u)$ .

The *computeLevel* function in Algorithm 4 goes over all nodes in gene tree  $G$  in the direction from the leaves to the root. The leaves of gene tree  $t \in L(G)$  are always mapped to its  $\sigma(t)$  and their variables are set according to it. For each inner node  $u \in V(G) \setminus L(G)$ , we recognize whether the node  $u$  has the same maximal mapping depth  $X[u]_{max}$  as  $\sigma(u)$  or the difference between  $X[u]_{max}$  and  $\sigma(u)$  is smaller than  $\epsilon$ , which allows us to determine if the node  $u$  is mapped to its  $\sigma(u)$  even when the depths are not same because of the rounding error. By default, the  $\epsilon$  is set to  $1 \times 10^{-6}$ .

In case that node  $u$  has not the same depth as  $\sigma(u)$  and also their *depthDifference* is bigger than  $\epsilon$ , we set *mappedToLca( $u$ )* to false. From the children of node  $u$ , we save the closest species node to the gene node  $u$  into variable *node* and run function *computeSpeciesNodeBelow* in Algorithm 2 to get *speciesNodeBelow( $u$ )*.

Otherwise, when the node  $u$  has the same depth as  $\sigma(u)$  or the *depthDifference* is smaller or equal to the  $\epsilon$ , we consider three cases can happen, where the first two cases check if both or at least one child is already mapped to  $\sigma(u)$ , because exactly one node from the gene tree can be mapped to exactly one node from species tree to be considerate as speciation. In the first case, we consider that both  $v \in children(u)$  have the same *speciesNodeBelow( $v$ )* as  $\sigma(u)$ . It means that both children of node  $u$

---

**Algorithm 4** Compute levels for nodes from gene tree  $G$ 


---

```

1: function COMPUTELEVEL( $u \in V(G)$ )
2:   if  $u \in L(G)$  then
3:     mappedToLca( $u$ ) = true
4:     speciesNodeBelow( $u$ ) =  $\sigma(u)$ 
5:      $l(u) = l(\sigma(u))$ 
6:   else
7:     for  $v \in \text{children}(u)$  do
8:       computeLevel( $v$ )
9:     depthDifference =  $D(\sigma(u)) - X[u]_{max}$ 
10:    if  $X[u]_{max} < D(\sigma(u))$  or depthDifference  $> \epsilon$  then
11:      mappedToLca( $u$ ) = false
12:      node =  $\max(v \in \text{children}(u) : \text{speciesNodeBelow}(v))$ 
13:      speciesNodeBelow( $u$ ) = computeSpeciesNodeBelow( $u$ , node))
14:    else
15:      if  $\forall v \in \text{children}(u) : \text{speciesNodeBelow}(v) = \sigma(u)$  then
16:        if  $\forall v \in \text{children}(u) : \text{mappedToLca}(v)$  then
17:           $v = v \in \text{children}(u) : \min(v)$ 
18:          levelDistanceFromChildren( $v$ )
19:          mappedToLca( $u$ ) = false
20:        else if  $\exists v \in \text{children}(u) : \text{speciesNodeBelow}(v) = \sigma(u)$  then
21:          mappedToLca( $u$ ) = false
22:        else
23:          mappedToLca( $u$ ) = true
24:          speciesNodeBelow( $u$ ) =  $\sigma(u)$ 
25:           $l(u) = l(\text{speciesNodeBelow}(u))$ 
26:          levelDistanceFromChildren( $u$ )

```

---

can be mapped to the same species node, so we check if their *mappedToLca*( $v$ ) is set to true. If the condition holds, we take the child  $v$  that is closer to the node  $u$  and set its *mappedToLca*( $v$ ) to false that also affects the distance from *children*( $v$ ), thus we call the *levelDistanceFromChildren*( $v$ ) to reset the distance. As the other child is already mapped to the species node, we set *mappedToLca*( $u$ ) to false.

The second case checks if at least one  $v \in \text{children}(u)$  is mapped to  $\sigma(u)$ . If the condition holds, we can not map another gene node to the species node, so we set *mappedToLca*( $u$ ) to false.

In the third case, any  $v \in \text{children}(u)$  has the same *speciesNodeBelow*( $v$ ) as the  $\sigma(u)$ , thus we set the *mappedToLca* to true.

Lastly, we set the *speciesNodeBelow*( $u$ ) to the  $\sigma(u)$  and level of node  $u$  to the level of computed *speciesNodeBelow*( $u$ ). Then, we run function *levelDistanceFromChildren* in Algorithm 3.

The *computeLevel* function goes over all nodes in the rooted gene tree  $G$ . If the gene tree  $G$  is balanced, the function *computeSpeciesNodeBelow* has running time  $O(\log N)$ , where  $N$  is number of nodes in gene tree  $G$ . The *levelDistanceFromChildren* function is computed in constant time. So the total preprocessing running time is  $O(N \log N)$  for balanced gene tree  $G$ . However, if the gene tree  $G$  is not balanced, the preprocessing running time is  $O(N^2)$ .

### 2.2.2 Main algorithm

The prerequisites for the *countDL* function shown in Algorithm 5 are computed in preprocessing. Besides, we need to have set the *countLossesAboveRoot* variable. If it is true, we count losses above root that occurred as a result of the gene tree  $G$  not containing genes of all species from the species tree  $S$ .

We count the evolutionary events in the direction from leaves to the root of gene tree  $G$ . For each node  $u \in V(G)$  and its parent  $v \in V(G)$ , we consider evolutionary events that occur in the node  $u$  and on the edge  $(u, v)$ . We do not compute evolutionary event in the parent, they are determined directly in the parent. In the root, we only calculate the evolutionary events that happened in the node as the root has no parent thus no edge to consider.

The number of duplications and gene losses are depicted as pair  $DL =$



**Algorithm 5** Counts duplications and gene losses in gene tree  $G$ 


---

```

1: function COUNTDL( $u \in V(G)$ )
2:   if  $u \neq L(G)$  then
3:      $v, w = \text{children}(u)$ 
4:      $DL_u = DL_v + DL_w$ 
5:      $\text{loss} = \text{levelDistanceFromParent}(u)$ 
6:     if  $\text{parent}(u) = \text{null}$  and  $\text{parent}(\sigma(u)) \neq \text{null}$  and  $\text{countLossesAboveRoot}$ 
       then
7:        $\text{loss} += l(u)$ ;
8:     if not  $\text{mappedToLca}(u)$  then
9:        $\text{duplication} = 1$ 
       return  $DL_u + (\text{duplication}, \text{loss})$ 

```

---

$(\text{duplication}, \text{loss})$ . The sum of the pairs  $DL_1$  and  $DL_2$  is computed as  $DL_1 + DL_2 = (\text{duplication}_1 + \text{duplication}_2, \text{loss}_1 + \text{loss}_2)$ . For each node  $u \in V(G)$ , we calculate the  $DL_u$ , which corresponds to the number of duplications and gene losses inferred in the subtree of node  $u$ . Thereafter, we infer the duplication and gene losses in the node  $u$  and on the edge above the node  $u$ .

The number of gene losses on the edge  $(u, v)$  is calculated as the number of species nodes on the path from node  $u$  to node  $v$ , which is precomputed in variable  $\text{levelDistanceFromParent}(u)$ . If  $\sigma(\text{root}(G)) \neq \text{root}(S)$ , thus the  $\sigma(\text{root}(G))$  is below  $\text{root}(S)$  means that some species do not have their gene in the gene tree. Therefore, the gene loss occurred before the  $\text{root}(G)$ , which resolves in extra gene loss that can be added if the variable  $\text{countLossesAboveRoot}$  is true. The truth value of  $\text{countLossesAboveRoot}$  is set by the user.

Duplication and speciation are easy to determine since it depends on whether the node  $u$  is mapped to  $\sigma(u)$  or above, which we already precomputed in  $\text{mappedToLca}(u)$ .

The function return pair  $DL$  corresponds to the number of duplication and gene losses in the subtree of node  $u$ , in the node  $u$  and on the edge above node  $u$ .

Gene losses and duplications are computed in constant time for one gene node. The  $\text{countDL}$  function goes over all nodes in gene tree  $G$  thus the running time is  $O(N)$ , where  $N$  is the number of all node in the gene tree  $G$ .

## 3 Implementation

The implementation of our software is built on Chládek's source code from his diploma thesis [5]. We apply his implementation of basic classes for defining the phylogenetic tree, parsing an unrooted gene tree and a rooted species tree from a file with several changes. We also use his implementation of the two-pass algorithm, which we mentioned before in Chapter 1.2.3.2, to compute the gene tree possible mapping depths.

On this basis, we implement our algorithms described in Chapter 2 to find the most parsimonious reconciliation.

Our software is implemented in the programming language Java and has a command-line interface.

### 3.1 Classes and variables

The implemented software consists of 15 classes that we briefly introduce:

- Class Node - represents a node in a phylogenetic tree
- Class Edge - represents an edge in a phylogenetic tree
- Class Interval - represents a length of edge as interval
- Class DL - represents a score of reconciliation
- Class UnrootedNode - represents a node in an unrooted gene tree
- Class UnrootedTree - represents an unrooted gene tree
- Class RootedExactNode - represents a node in a rooted species tree with exact branch lengths

- Class `RootedExactTree` - represents a rooted species tree with exact branch lengths
- Class `RootedIntervalNode` - represents a node in a rooted gene tree with inexact branch lengths
- Class `RootedIntervalTree` - represents a rooted gene tree with inexact branch lengths
- Class `Parser` - parses phylogenetic trees and their leaf-mapping from files
- Class `Loader` - load arguments from a command line and calls `Parser` on files with stored phylogenetic trees
- Class `Printer` - prints reconciliation solutions into files
- Class `Reconciliator` - computes reconciliation with its score
- Class `Main` - the main class of the software that calls other classes to load files, compute and print reconciliation

## 3.2 Differences from the original source code

As our source code is built on Chládek's source code, we show the differences in relations and entities between source codes in Fig. 3.2. The differences, new variables and classes are described below by classes in more details.

### **Class `Interval`**

We change the name of parameters to *minLength* and *maxLength* and delete the unnecessary variable *OriginalMappingDepth* as we use the *Interval* to store the minimal and maximal length of new possible edges that can be created after subdividing the original root edge in function *getIntervals* (Algorithm 1) in *Reconciliator* class. The parameters are numbers of type double.

### **Class `RootedExactNode`**

In the *RootedExactNode* class, we define a new *level* variable. It is a number of type integer. The variable is set while parsing the species tree in the *Parser* class and used for computing the level of gene nodes as shown in function *computeLevel* (Algorithm

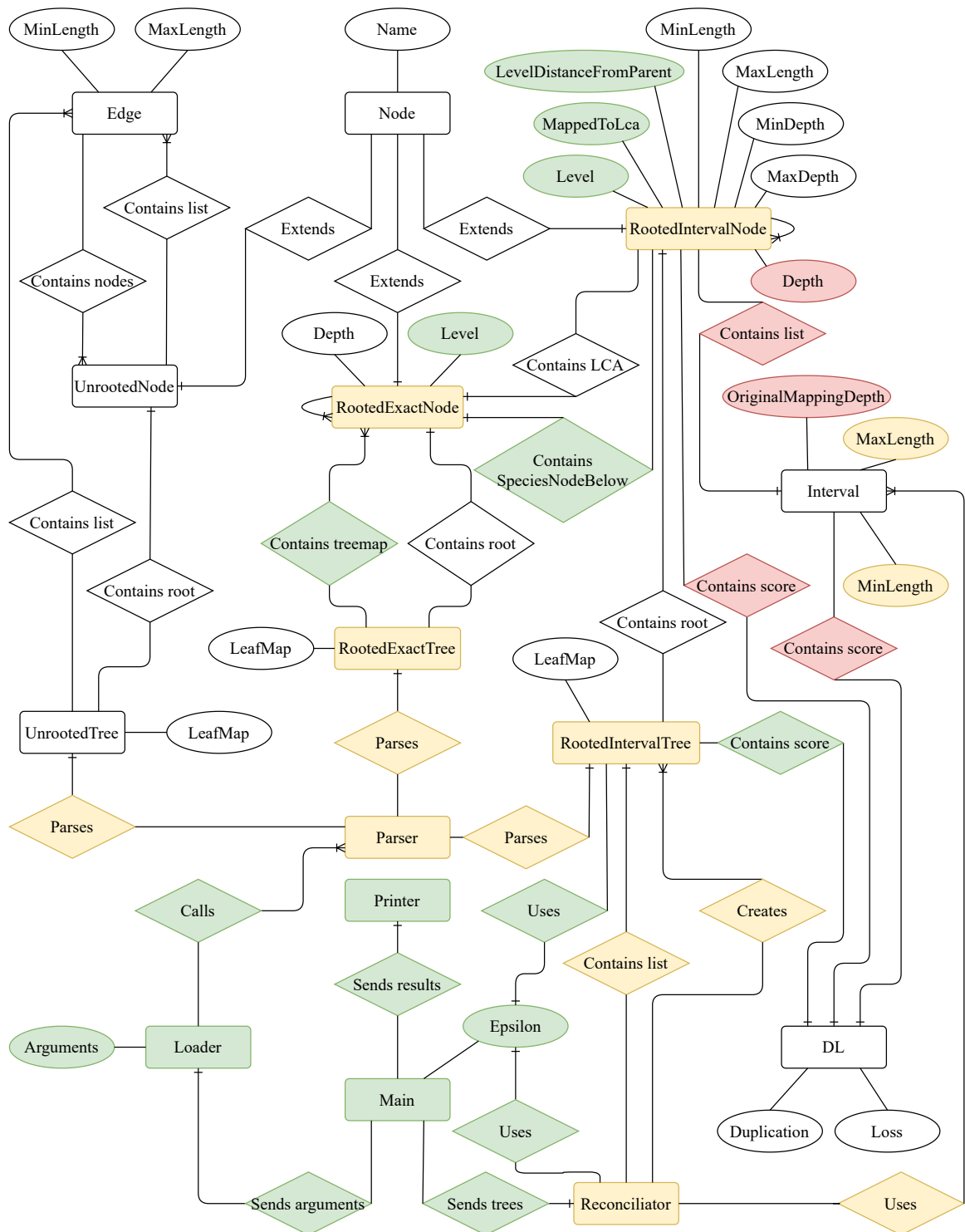


Figure 3.1: Red entities and relations represent deleted variables or functions. Yellow entities and relations show changes to the original source code. Green entities and relations depict added classes, variables or functions.

4) in *Reconciliator* class.

## Class RootedExactTree

We add structure *TreeMap* containing leaves, where the key is the name of the leaf and the value is an object representation of that leaf as *RootedExactNode*. It is used to set the LCA-mapping variable of leaves in a gene tree according to the leaf-mapping.

### **Class RootedIntervalNode**

The unnecessary *depth* variable was deleted as we store the mapping depth of a *RootedIntervalNode* in variables *minDepth* and *maxDepth*. We also delete the *DL*, because we do not store the reconciliation score in nodes and the list of *Interval* since we have the mapping depth interval and the interval of the edge above stored in separated variables. Besides, we add new variables: *level*, *mappedToLca*, *levelDistanceFromParent* and *speciesNodeBelow* that are computed and used in algorithms in Chapter 2.2. The *level* and *levelDistanceFromParent* are numbers of type integer. The *mappedToLca* is data type boolean *speciesNodeBelow* is object of type *RootedExactNode*.

### **Class RootedIntervalTree**

In the *RootedIntervalTree*, we implement our algorithms from Chapter 2.2 and their are called from the *Reconciliator* class. We add the variable of type *DL* to store the number of duplication and gene losses inferred in the gene tree.

### **Class Parser**

We change the parsing method for a species tree and an unrooted gene tree. In the species tree parsing method, we infer *level* in species tree nodes. In the unrooted gene tree parsing method, we resolve inexact branch lengths if the given gene tree has inexact branches or we transform exact branch lengths to inexact branch lengths if the given gene tree has exact branch lengths. To modify the exact branch lengths into inexact branch lengths, the user needs to set the tolerance value that has to be from interval  $\langle 0, 1 \rangle$ . The interval of the edge is compute as:  $\langle length - (tolerance \cdot length), length + (tolerance \cdot length) \rangle$ .

Furthermore, we add a parsing method for a rooted gene tree that can either resolve inexact branch lengths or transform exact branch lengths to inexact branch lengths. We included a method to parse mapping of gene leaves to species leaves from a file and a method for transforming a rooted gene tree to an unrooted gene tree. This allows us to forget the original root of a gene tree and find a new one with function *getIntervals* (Algorithm 1). All methods in the class are called from the *Loader* class.

### **Class Loader**

The *Loader* class is initialized with arguments from the command line. It loads the arguments and calls functions from *Parser* according to the requirements specified in the arguments. We recognize 11 arguments that are listed in Table 3.1.

Table 3.1: Input arguments

Argument	Description
-help	shows help information with all possible input arguments
-S <species tree>	path to the rooted species tree file in Newick tree format (mandatory input)
-G <gene tree>	path to the rooted or unrooted gene tree in Newick tree format (mandatory input)
-M <species map>	mapping of genes to species
-t <tolerance>	required tolerance from interval $\langle 0, 1 \rangle$ (By default, it is set to 0.5)
-s <step>	required step (By default, it is set to 0.5)
-r	signifies that the given gene tree is rooted (By default, the given gene tree is considered to be unrooted)
-reroot	signifies that the given rooted gene tree is wished to be rerooted
-l	signifies counting the gene losses above the root of given gene tree
-p <print type>	required print type from two options " <i>sol</i> " and " <i>rel</i> " (By default, the print type is set to be " <i>sol</i> ")
-epsilon <epsilon>	required epsilon as number of type double (By default, it is set to $1 \times 10^{-6}$ )

The main method is *loadArgs* which is called from the *Main* class and returns an array of *Object*: a rooted species tree, a rooted gene tree or an unrooted gene tree, *step* and *countLossesAboveRoot* needed for algorithms from Chapter 2.2 implemented in *Reconciliator* class along with *dirPathGene* and *printType* required in *Printer* class functions.

### Class Reconciliator

The algorithms from Chapter 2 and the two-pass algorithm by Chládek from Chapter 1.2.3.2 are called in the *Reconciliator* class. These are the main methods to obtain the most parsimonious reconciliation.

The *Reconciliator* class is initialized with a rooted species tree, a rooted gene tree or an unrooted gene tree, *step* and *countLossesAboveRoot* that are inferred in *Loader* class. It uses *Interval* to store the minimal and maximal length of new edges in function *getIntervals* (Algorithm 1), which are used for subdividing the root edge and creating the *RootedIntervalTree*. The class returns a list of *RootedIntervalTree* with most par-

simonious reconciliation that has its reconciliation score store in *DL*.

### **Class Printer**

The class is initialized with list of *RootedIntervalTree*, *dirPathGene* and *printType*. We recognise two types of printing the results of the reconciliation into the file: "*rel*" and "*sol*".

With the "*rel*" type, the function prints relations between genes in the gene tree for each inferred *RootedIntervalTree*. If solutions contain more trees with the same topology and reconciliation score, it prints the number of the same tree after each gene relation printout. For each  $u \in L(G)$ , it prints "*gene*" and name of the node  $u$ . For each  $v \in V(G) \setminus L(G)$ , it allows evolutionary events as speciation, duplication and gene loss depending on what occurred in the node  $v$  or on the branch from the node  $v$  to its parent  $parent(v)$ . The speciation and duplication are inferred in the node  $v$ . They print as "*spec*" for speciation and "*dup*" for duplication with names of  $children(v)$  separated by a tab. The gene loss happens on the branches. It prints as "*loss*" with the name of node  $v$ . The name of each node consists of genes found in the subtree of that node separated by commas.

The "*sol*" type prints gene trees in special Newick format, where branch lengths are inexact. All computed solutions are printed into one *.txt* file, where solutions are separated with an empty line.

The file with results is saved in the same directory as the given gene tree.

### **Class Main**

The *Main* class receives arguments given by user and send it to the *Loader* class for processing. Sequentially, sends a rooted species tree, a rooted gene tree or an unrooted gene tree, *step* and *countLossesAboveRoot* to the *Reconciliator* class to infer the most parsimonious reconciliation, which returns list of *RootedIntervalTree*. Eventually, it sends the list of *RootedIntervalTree*, *printType* and *dirPathGene* to the *Printer* to print the results. Also, it stores the value of  $\epsilon$  that are required in functions *getIntervals* (Algorithm 1) in *Reconciliator* class and *computeLevel* (Algorithm 4) in *RootedIntervalTree* class.

## 4 Experiments

We evaluate our software for isometric reconciliation on a simulated and real dataset that were used to evaluate SPIMAP [17] and TreeFix [20] in previous studies. The simulated dataset consists of 1000 simulated gene families of two clades of species: 12 *Drosophila* genomes and 16 fungal genomes, generated with the SPIMAP model. The real dataset includes 5351 real gene families from 16 fungal genomes.

To compare, we run other software for computing the gene tree reconciliation: Notung [8], TreeFix [20] and Treerecs [7], on the same dataset.

### 4.1 Simulated dataset

In the simulated dataset, we know the correct gene tree with its evolutionary events, which allows us to test several aspects. In the experiments, we measure the accuracy of correctly inferred topology, branches, duplications and gene losses. For each software, we measure the time of computing the results.

We evaluate our software for isometric reconciliation for two cases: a rooted gene tree with different tolerance settings and an unrooted gene tree with different step and tolerance setting. In the first case, we take the rooted gene tree as it is from the dataset and run a reconciliation with *countDL* function (Algorithm 5 in Chapter 2.2.2). We evaluate it several times with different tolerance setting to scale up the edges. The tolerance setting are 0.0, 0.000001, 0.00001, 0.001, 0.1, 0.3, 0.5, 1.0 and the  $\epsilon = 1 \times 10^{-6}$ .

The second case takes the rooted gene tree as input with an argument to reroot the given gene tree. We forget the root of the rooted gene tree and transform it into the unrooted gene tree. On the obtained unrooted gene tree, we run the *getIntervals* function (Algorithm 1 in Chapter 2.1) to subdivide each edge of the unrooted gene tree and then run a reconciliation with *countDL* function (Algorithm 5 in Chapter 2.2.2).



We want to find a new root minimizing the number of inferred duplications and gene losses in reconciliation. The process is executed several times with different tolerance and step settings. The tolerance and  $\epsilon$  setting are the same as in the first case. The step setting are 0.0, 0.3, 0.5, 1.0, 2.0.

## 4.2 Real dataset

The correct gene tree is unknown in the real dataset, thus we use different metrics as with the simulated dataset. We measure the number of inferred duplications and gene losses overall gene trees in the dataset and the consistency score of duplications.

The duplication consistency score  $dcs$  shows the plausibility of inferred duplications. For a duplication node  $u$  with children  $v$  and  $w$ , the duplication consistency score  $dcs(u) = (L \cup R) \mid (L \cap R)$  is computed as the union of  $L$  and  $R$  over intersection of  $L$  and  $R$ , where  $L$  is the set of species represented in left child  $v$  and  $R$  is the set of species represented in right child  $w$ .

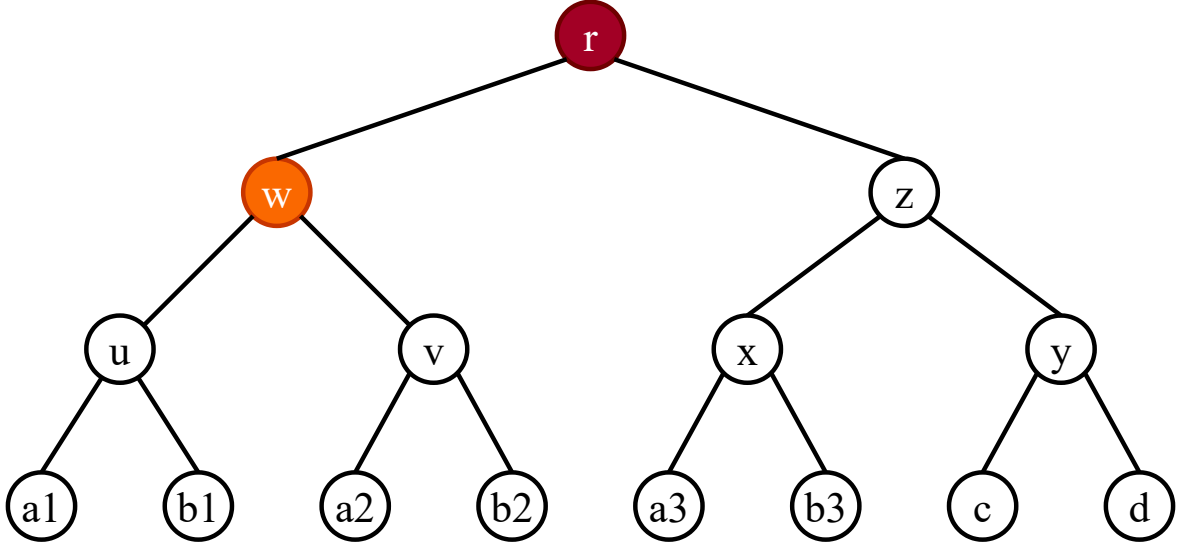


Figure 4.1: **Duplication consistency score for node  $w$**

The set of species of left child  $u$  is  $L = \{A, B\}$  and the set of species of right child  $v$  is  $R = \{A, B\}$ . The duplication consistency score is computed as:  $dcs(w) = (L \cup R) \mid (L \cap R) = 2 \div 2 = 1$ .

**Duplication consistency score for node  $r$**

The set of species of left child  $w$  is  $L = \{A, B\}$  and the set of species of right child  $z$  is  $R = \{A, B, C, D\}$ . The duplication consistency score is computed as:  $dcs(r) = (L \cup R) \mid (L \cap R) = 2 \div 4 = 0.5$ .

## 5 Results

# Conclusion

# Appendix

# Bibliography

- [1] Mukul Bansal. Tutorial: Treefix and treefix-dtl. Available from: <http://compbio.mit.edu/treefix/tutorial.html>, 2014. [Accessed 12 Jan 2021].
- [2] Bastien Boussau. Phyllog: joint reconstruction of species and gene phylogenies. Available from: <https://pbil.univ-lyon1.fr/software/phyllog/>, 2012. [Accessed 12 Jan 2021].
- [3] Bastien Boussau et al. Genome-scale coestimation of species and gene trees. *Genome Research*, 23(2):323–330, Feb. 2013.
- [4] Broňa Brejová et al. Isometric gene tree reconciliation revisited. *Algorithms for Molecular Biology*, 12(1):17, Jun. 2017.
- [5] Radoslav Chládek. Algorithms for isometric gene tree reconciliation. Master’s thesis, Comenius University in Bratislava, 2019.
- [6] Radoslav Chládek et al. Isometric gene tree reconciliation with software.interval software.edge lengths.
- [7] Nicolas Comte et al. Treerecs: an integrated phylogenetic tool, from sequences to reconciliations. *bioRxiv*, Oct. 2019.
- [8] Dave Danicic et al. *Notung 2.8: A Manual*. Notung Development Team, Mar. 2015.
- [9] Jean-Philippe Doyon, Cedric Chauve, and Sylvie Hamel. Algorithms for exploring the space of gene tree/species tree reconciliations. In *Nelson C.E., Vialette S. (eds) Comparative Genomics. RECOMB-CG 2008. Lecture Notes in Computer Science*, volume 5267, pages 1–13. Springer-Verlag, Berlin, Heidelberg, Oct. 2008.

- [10] Joseph Felsenstein. *Inferring Phylogenies*. Sinauer, Sunderland, 2003.
- [11] Morris Goodman et al. Fitting the gene lineage into its species lineage, a parsimony strategy illustrated by cladograms constructed from globin sequences. *Systematic Zoology*, 28(2):132–163, Jun. 1979.
- [12] Damir Hasić and Eric Tannier. Gene tree species tree reconciliation with gene conversion. *Journal of Mathematical Biology*, 78(6):1981–2014, May 2019.
- [13] Li Heng. Treesoft: Treebest. [online]. Available from: <http://treesoft.sourceforge.net/treebest.shtml>. [Accessed 12 Jan 2021].
- [14] INRIA. Tutorial – treerecs. Available from: <https://project.inria.fr/treerecs/tutorial/>, 2011. [Accessed 12 Jan 2021].
- [15] Jian Ma et al. The infinite sites model of genome evolution. *Proceedings of the National Academy of Science*, 105(38):14254–14261, Sep. 2008.
- [16] Matthew D. Rasmussen. Spimap documentation. Available from: <http://compbio.mit.edu/spimap/pub/spimap/doc/spimap-manual.html>, 2011. [Accessed 12 Jan 2021].
- [17] Matthew D. Rasmussen and Manolis Kellis. A bayesian approach for fast and accurate gene tree reconstruction. *Molecular Biology and Evolution*, 28(1):273–290, Jan. 2011.
- [18] Benjamin Vernot et al. Reconciliation with non-binary species trees. *Journal of Computational Biology*, 15(8):981–1006, Oct. 2008.
- [19] Albert J. Vilella et al. Ensemblcompara genetrees: Complete, duplication-aware phylogenetic trees in vertebrates. *Genome Research*, 19(2):327–335, Feb. 2009.
- [20] Wu Yi-Chieh et al. Treefix: Statistically informed gene tree error correction using species trees. *Systematic Biology*, 62(1):110–120, Jan. 2013.
- [21] Wu Yi-Chieh et al. Treefix. Available from: <https://www.cs.hmc.edu/~yjiw/software/treefix/>, 2014. [Accessed 12 Jan 2021].