

# PHYLOGENETIC SUPERTREE CONSTRUCTION USING TRIPLETS

## A DISSERTATION

*Submitted in partial fulfillment of the  
requirements for the award of the degree*

*of*

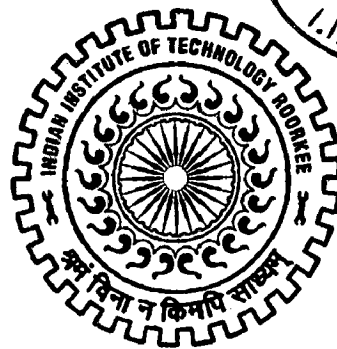
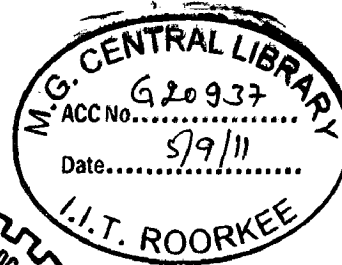
MASTER OF TECHNOLOGY

*in*

COMPUTER SCIENCE AND ENGINEERING

*By*

**SATYA GOVINDU AMUJURU**



DEPARTMENT OF ELECTRONICS AND COMPUTER ENGINEERING  
INDIAN INSTITUTE OF TECHNOLOGY ROORKEE  
ROORKEE - 247 667 (INDIA)

JUNE, 2011

## CANDIDATE'S DECLARATION

---

I hereby declare that the work, which is being presented in the dissertation entitled “**PHYLOGENETIC SUPERTREE CONSTRUCTION USING TRIPLETS**” towards the partial fulfillment of the requirement for the award of the degree of **Master of Technology in Computer Science and Engineering** submitted in the Department of Electronics and Computer Engineering, Indian Institute of Technology Roorkee, Roorkee, Uttarakhand (India) is an authentic record of my own work carried out during the period from July 2010 to June 2011, under the guidance of **Dr. Rajdeep Niyogi, Assistant Professor**, Department of Electronics and Computer Engineering, IIT Roorkee.

The matter presented in this dissertation has not been submitted by me for the award of any other degree of this or any other Institute.

Date: 29/06/2011

Place: Roorkee

  
(SATYA GOVINDU AMUJURU)

---

## CERTIFICATE

This is to certify that the above statement made by the candidate is correct to the best of my knowledge and belief.

Date: 29/06/2011

Place: Roorkee

  
(Dr. Rajdeep Niyogi)

Assistant Professor

Department of Electronics and Computer Engineering

IIT Roorkee.

## **ACKNOWLEDGEMENTS**

---

First and foremost, I would like to extend my heartfelt gratitude to my guide and mentor **Dr. Rajdeep Niyogi**, Professor, Department of Electronics and Computer Engineering, Indian Institute of Technology Roorkee, for his invaluable advices, guidance, encouragement and for sharing his broad knowledge. His wisdom, knowledge and commitment to the highest standards inspired and motivated me. He has been very generous in providing the necessary resources to carry out my research. He is an inspiring teacher, a great advisor, and most importantly a nice person.

I am greatly indebted to all my friends, who have graciously applied themselves to the task of helping me with ample moral supports and valuable suggestions.

On a personal note, I owe everything to the Almighty and my parents. The support which I enjoyed from my father, mother and other family members provided me the mental support I needed.

**SATYA GOVINDU AMUJURU**

## ABSTRACT

---

Phylogenetics deals with evolutionary relatedness between various species. Origin of species can be established by studying this relationship. The central task in phylogenetics is to infer this relationship among a given set of species. These relationships are usually represented by a phylogenetic tree with the species of interest at the leaves and where the internal vertices of the tree represent ancestral species. The amount of available molecular data is increasing exponentially and, given the continual advances in sequencing techniques and throughput, this explosive growth will likely to continue.

Biologists assemble large multi-gene data set from available vast amount of data for use in phylogenetic analyses which imposes distinct computational challenges. As this biological data is vast it is infeasible to construct a large tree for analysis. As a result, supertree methods have been the focus of much research. Supertree methods comprise one approach to reconstructing large phylogenies, given estimated trees for overlapping subsets of the entire set of taxa, using various algorithmic techniques.

Several supertree methods have previously been developed. In this report implementation of Triplet Inference and Local Inconsistency (TILI) and Triplet Supertree Heuristic algorithms has been done. Comparative study is carried out on the basis of execution time and Accuracy measurement. The obtained results infer that Triplet Supertree Heuristic algorithm is faster and accurate than the TILI.

# Table of Contents

<b>Candidate’s Declaration &amp; Certificate.....</b>	<b>i</b>
<b>Acknowledgements.....</b>	<b>ii</b>
<b>Abstract.....</b>	<b>iii</b>
<b>Table of Contents.....</b>	<b>iv</b>
<b>List of Figures.....</b>	<b>vi</b>
<b>List of Tables.....</b>	<b>vii</b>
<b>1. Introduction</b>	<b>1</b>
1.1 Introduction of Phylogenetics.....	1
1.2 Motivation.....	4
1.3 Statement of the Problem.....	5
1.4 Organization of the Report.....	5
<b>2. Background and Literature Review</b>	<b>6</b>
2.1 Basic Definitions .....	6
2.1.1 Trees.....	6
2.1.2 Subtrees.....	7
2.1.3 Clusters and Bipartitions.....	8
2.1.4 Triplets and Quartets.....	10
2.2 Phylogenetic trees.....	11
2.3 Phylogenetic supertrees.....	12
2.4 The process of a Phylogenetic Study.....	14
2.5 Literature Review.....	16
<b>3. Phylogenetic Supertree Construction using Triplet (Inference and ) Local Inconsistency</b>	<b>20</b>
3.1 Triplet Terminology.....	20

3.2	T(I)LI Algorithm.....	21
3.3	Implementation and Results.....	26
<b>4.</b>	<b>Heuristic approach for Phylogenetic Supertree</b>	<b>29</b>
4.1	Preliminaries.....	30
4.2	Tree Editing Operations.....	31
4.2.1	ReRoot operation.....	31
4.2.2	Tree Bisection and Reconnection .....	32
4.2.3	Subtree Pruning and Regrafting .....	33
4.3	Triplet Supertree Heuristic Algorithm.....	33
4.4	Implementation and Results.....	36
<b>5.</b>	<b>Conclusion and Future Work</b>	<b>39</b>
5.1	Conclusion.....	39
5.2	Future Work.....	39
	<b>REFERENCES.....</b>	<b>40</b>

## LIST OF FIGURES

---

Figure No.	Figure Name	Page No.
Figure 2.1	An example of rooted and unrooted trees .....	07
Figure 2.2	An Example of a subtree.....	08
Figure 2.3	The four rooted trees with three leaves .....	10
Figure 2.4	The four unrooted trees with four unrooted leaves .....	10
Figure 2.5	Rooted Phylogenetic Tree.....	11
Figure 2.6	An example of input source trees.....	12
Figure 2.7	A possible supertree for the set of input source trees .....	13
Figure 3.1	Triplet with four non trivial topologies.....	20
Figure 3.2	Algorithm for Decomposing forest into Triplets .....	21
Figure 3.3	Algorithm for inferring missed triplets.....	23
Figure 3.4	An extension of Triplet to a quartet by adding a leaf .....	24
Figure 3.5	A triplet that has been expanded displays three other triplets .....	24
Figure 3.6	Algorithm for composing triplets to a supertree .....	25
Figure 3.7	The connecting paths between the leaf nodes of a resolved triplet in a tree.....	26
Figure 3.8	Runtime comparison between TLI using both DM and LCA and TILI .....	28
Figure 4.1	Reroot operation.....	31
Figure 4.2	Tree Bisection and Reconnect operation.....	32
Figure 4.3	Algorithm for SPR-S problem .....	34
Figure 4.4	Algorithms for Preprocess TripletSum and Extended TripletSum .....	35
Figure 4.5	Runtime comparisons between TILI and TH .....	37

## LIST OF TABLES

---

Table No.	Table Name	Page.No
Table 3.1	Time in seconds required to build supertrees from forests of varying sizes.....	27
Table 4.1	Time comparison between TILI and TH for constructing supertree.....	36
Table 4.2	Triplet Similarity Score between TILI and TH.....	37



# Chapter 1

## Introduction

---

### 1.1 Introduction of Phylogenetics

A fundamental concept of the theory of evolution, independently developed by Charles Robert Darwin and Alfred Russel Wallace is that species share a common origin and have subsequently diverged through time. Interestingly, both men came to use the metaphor of a great tree to illustrate this notion of descent with modification; ever since, biologists have been using treelike diagrams to describe the pattern and timing of events that gave rise to the earth's biodiversity.

According to Pennisi today biologists have catalogued about 1.7 million of species and they estimates of the total number of species ranges from 4 to 100 million [1]. With this explosion in the amount of data in taxonomy it is no longer possible to analyze and build trees by hand. The fact that all life on Earth is genetically related is one of the most profound scientific observations of all time.

Scientists have observed that genealogical relationships among living things can (generally) be represented by a vast evolutionary tree. Constructing this tree of all life is a fundamental scientific problem facing human kind today. Tree Of Life (TOL) is a project (refer <http://www.tolweb.org/tree>) to discover the relationships among all the species on Earth. This project presents one of the greatest challenges of science in reconstructing the evolutionary history of every living organism on the earth.

Under TOL many smaller projects such as “The Green Tree of Life”, and “Assembling the Fungal Tree of Life (AFTOL)” are also proposed. Huge support is also provided from different funding agencies and countries to the development of computational methods to reconstruct the tree of life (refer <http://www.phylo.org>). The field responsible for this undertaking is phylogenetics. Trees that depict evolutionary relationships between species, or other entities, are called phylogenetic trees or phylogenies.

In earlier days, phylogenies were built using morphological data i.e. for classifying seabirds might involve comparison of beak shapes or other distinct physical attributes. With the discovery

of DNA and the design of sequencing techniques, a new kind of information became available: molecular data. The increase in the availability of DNA and protein sequence data has increased interest in molecular phylogenetics and classification. Molecular phylogenetics overcomes limitations of morphological phylogenetics such as, convergent evolution, finding the relationship among bacteria and the comparison of distinctly related organisms [2].

The successes of molecular biology and genetics have yielded large amounts of data that are increasingly quantitative in nature, both in number of taxa as well as sequence length, available for phylogenetic construction. Most of the inferences in comparative biology depend on accurate estimates of evolutionary relationships [3, 4]. Consequently there is a growing need for new techniques to speed up this process accurately.

Reconstruction of ancestral relationships from contemporary data is widely used to provide evolutionary and functional insights into biological system. These insights are largely responsible for the development of new crops in agriculture, drug design and in understanding the ancestors of different species.

The other applications of phylogenetics include:

- In tracking the origins and development of humans over time [5]
- Understanding of the process of evolution of different species, constraints, behavior and evolutionary time [6, 7].
- Performing functional prediction of genes. Most of the similarly looking genes show the similar functionality, looking at the evolutionary history of the newly sequenced gene will help in predicting its functionality [8].
- Development of vaccines, antimicrobial and herbicides [9, 10].
- Forensics studies [11].

In a typical molecular phylogenetic analysis, a tree on a particular set of species is constructed by first collecting the DNA or protein sequences for a homologous gene in each species, and then using those sequences to construct a tree on that set of species. (Homology is shared evolutionary history, and genes are called homologous if they descended from a common ancestor gene.) Using this sort of process to construct the Tree of Life is not feasible.

There are two main approaches for analyzing large number of gene datasets: constructing a tree by analyzing the entire dataset as a whole, or constructing a tree on each gene dataset separately and then combining those trees into a single tree on the entire set of species. The first one is called the combined analysis, “total evidence” approach, and the second one is called the supertree approach, and a method that combines trees with overlapping leaf sets into a single tree is called a supertree method.

Supertrees are meant to contain the relations between all species in a set of smaller trees, which are combined using a supertree method. There are several advantages to supertree methods. The first is that whenever summarizing the results of available studies on (subset of) a particular group of interest when access to the sequences is not possible. Second, producing a tree from disparate data-types, such as molecular, morphological, and gene-order data, that requires independent types of analysis (and therefore prohibits a combined analysis). But supertrees can be created from those small trees using any one of the supertree method. Since how the input trees have been made is not important to a supertree method. Third is that if it is infeasible to infer a tree for a large set of species in one go, it might be possible to infer trees for subsets of the species. These can then be combined in a supertree in manageable time.

The use of supertree methods to solve above such problems has received much criticism [12]. For example, many supertree methods have been shown to have a size bias, favoring the relationships supported by larger input trees over smaller ones. Some have also been shown to have a shape bias, producing either balanced or unbalanced trees more often [13], and the input trees on which supertrees are constructed may not be completely independent (some primary data may have been used in the inference of more than one input tree). However, given the limitations of current sequence-based phylogenetic reconstruction methods, supertrees are presently a necessary tool for many phylogenetic problems [14]. In fact, most, if not all, detailed estimates of the Tree of Life to date has been constructed using a supertree approach: using informal supertree methods, splicing together various phylogenetic trees and tree representations of taxonomies, to achieve a single tree on all currently known species.

## 1.2 Motivation

The central task in phylogenetics is to infer the evolutionary relationships among a given set of species. These relationships are usually represented by a phylogenetic tree with the species of interest at the leaves and where the internal vertices of the tree represent ancestral species.

The large applications and huge available sequence data, reflect the need for developing fast and accurate methods for constructing the phylogenies. Most of the existing methods accept molecular sequence data or distance data for the construction of the phylogenetic history.

Advances in molecular biology and genomics lead to the large amounts of data, both in number of taxa as well as sequence length, available for phylogenetic construction. Moreover, most of the inferences in comparative biology depend on accurate estimates of evolutionary relationships. In this scenario, methods, such as maximum likelihood, parsimony search, or quartet puzzling, play an important role. But they are very slow and cannot scale up to the size of the genetic data available. These methods compare trees according to a specific criterion. Ideally, all possible trees should be compared. But they are very slow and cannot scale up to the size of the genetic data available.

After generating the genomic data of thousands of living organisms, now comes the task of classifying them and making a single tree of life. This is a challenging task and till date no algorithm exists that can compute most accurate tree of life. The existing phylogenetic reconstruction algorithms cannot be used for this task as they suffer with poor efficiency and computational hardness problems.

The problem with the phylogenetic tree construction approach is that they cannot be used for the classification of 1.7 million described species on the earth. The existing methods cannot be used for constructing the tree of life due to aforementioned shortcomings. However, given the limitations of current phylogenetic reconstruction methods, supertrees are presently a necessary tool for such phylogenetic problems. In fact, most, if not all, detailed estimates of the Tree of Life can only be done using phylogenetic supertree methods.

### **1.3 Problem statement**

The aim of this thesis is to implement the construction of phylogenetic tree using triplet tree algorithms.

1. Implementation of Phylogenetic Supertree using Triplet (Inference and) Local Inconsistency.
2. Implementation of phylogenetic supertree using Triplet supertree heuristics.

### **1.4 Organization of the Report**

The report comprises of five chapters including this chapter introduces the topic and states the problem. The rest of the dissertation report organised as follows

Chapter 2 gives an overview of basic trees, phylogenetic trees and supertrees and the process of phylogenetic study. Finally it gives a brief literature review of the related work.

Chapter 3 gives details about terminology of the triplets. It discusses about brief overview and implementation of the triplet inference and local inconsistency algorithm.

Chapter 4 gives details about hill climbing heuristics, operations that can be done on the trees. It discusses about approach and implementation of Algorithm.

Chapter 5 concludes the dissertation work and gives suggestions for future work.

## Chapter 2

### Background and Literature Review

---

The biological discipline dedicated to reconstructing organismal phylogenies is called phylogenetics. The term phylogenetics derives from the Greek, *phylon* means race or old family, and *genos* means birth or origin, *gennetikos* means related to generation or the genesis of something. The branching pattern of the tree represents the splitting of biological lineages, and the lengths of the branches can be used to signify the age of those events. Today, biologists call these treelike diagrams phylogenies. Reconstruction of these phylogenies can provide an evolutionary framework for studying a variety of problems led to their application in almost every other sub discipline of biology. Since evolutionary and this thesis is about phylogenetic trees, it is therefore appropriate to start by defining from a tree.

#### 2.1 Basic Definitions

##### 2.1.1 Trees

A tree consists of internal nodes and leaf nodes, and they are connected by edges which can have edge lengths. A tree can be a binary or non-binary tree. One node may be designated as a root node, in which case the tree is rooted otherwise it is unrooted.

A binary rooted tree can be treated as a directed acyclic connected graph having exactly one internal vertex of degree two and remaining internal vertices are having degree more than two. The leaves are vertices of degree one. In rooted trees there are directed edges which are directed away from the root. All nodes except the root have a parent, and all nodes except the leaf nodes have a number of child nodes. The nodes on the path from a node  $x$  to the root are called ancestors of  $x$ . A descendant of  $x$  is any node which has  $x$  as an ancestor.

A binary unrooted tree is a tree with no special node with degree two and it can be treated as an acyclic connected graph having no internal vertices of degree two and leaves are vertices of degree one. An unrooted tree may be rooted by selecting any node as the root. The set of nodes directly connected to a node  $x$  by edges in any direction is called the adjacent of  $x$ .

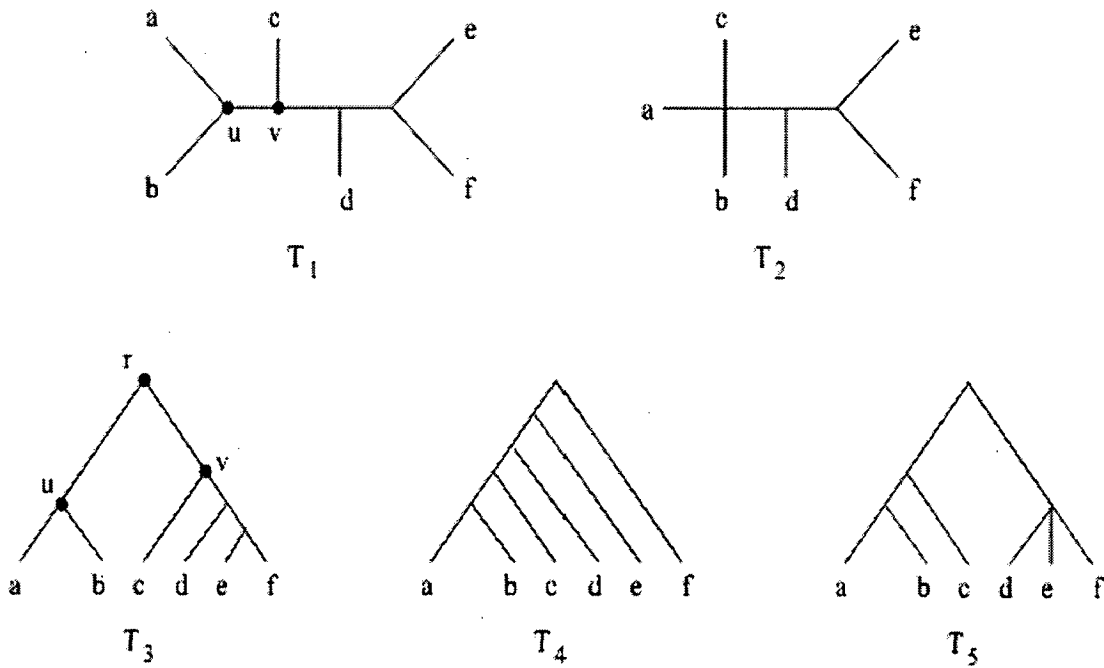


Figure 2.1: An example of rooted and unrooted trees

In Figure 2.1,  $T_1$  and  $T_2$  are unrooted,  $T_3$ ,  $T_4$ , and  $T_5$  are rooted,  $T_1$ ,  $T_3$ , and  $T_4$  are binary, and  $T_2$  and  $T_5$  are non-binary. Notice that  $T_1$  can be obtained from  $T_3$  or  $T_4$  by suppressing the root vertex of either of these rooted trees. Equivalently,  $T_3$  and  $T_4$  can be obtained from  $T_1$  by subdividing an edge in  $T_1$  with a root vertex. For tree  $T_3$  internal nodes are  $\{r, u, v\}$ , leaves are  $\{a, b, c, d, e, f\}$ , and adjacent nodes of  $v$  are  $\{r, c, d, e, f\}$ .

In a fully resolved binary tree with  $n$  leaf nodes there are  $n-1$  internal nodes. In Figure 2.1,  $T_3$  and  $T_4$  are fully resolved binary trees. The topology of a tree is the arrangement of the nodes and edges, ignoring edge lengths and symmetrical differences (i.e. it is arbitrary in which order the children of a node appear). If all leaf nodes are connected to a single internal node, then the tree is called a star tree because of its shape.

### 2.1.2 Subtrees

Let  $T$  be a rooted tree and choose a vertex  $v$  in  $T$ . If we remove the edge between  $v$  and the parent of  $v$ , say  $p$ , we get two connected subgraphs. Then let  $v$  be the root of the subgraphs containing  $v$ , then this is called the subtree of  $T$  rooted at  $v$ . Briefly a subtree  $T'$  is a tree whose

vertices and edges form the subsets of the vertices and edges of a given tree  $T$ . An example of a subtree is shown in Figure 2.3 in which tree  $T_1$  is a subtree of tree  $T$ .

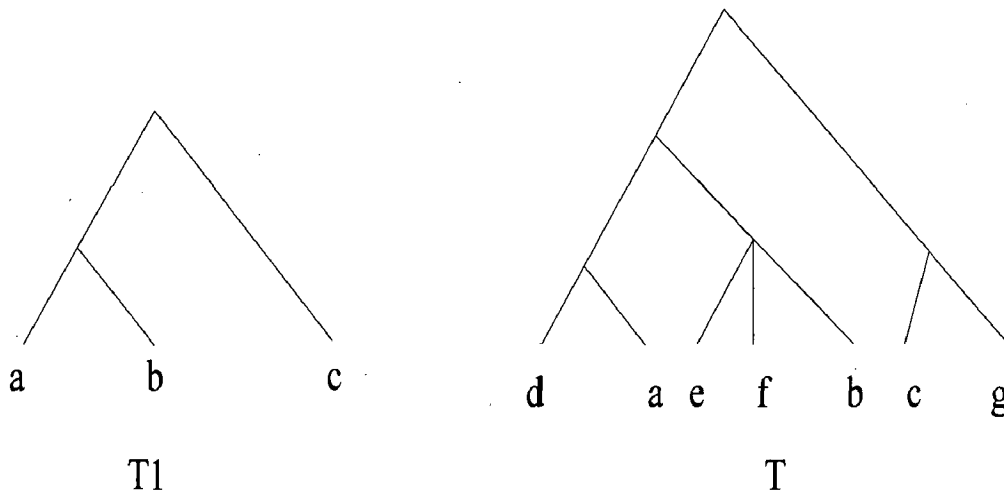


Figure 2.2: An Example of a subtree

### 2.1.3 Clusters and Bipartitions

Each vertex in a rooted tree defines a group of taxa that are more closely related to each other than they are to any other taxon in the tree; such a group is called a clade. The clade defined by a vertex  $v$  is the set  $C(v)$  containing all the taxa that can be reached from the root via a path containing the vertex  $v$ .

In phylogeny, such a group is often called a cluster. The set of all clusters defined by some vertex in a tree  $T$  is denoted  $H(T)$ ; i.e.,  $H(T) = \{C(v) : v \in V(T)\}$ . Notice that  $H(T)$  is a hierarchy, which means that for every two members  $C_1, C_2$  of  $H(T)$ , either  $C_1 \subseteq C_2$ ,  $C_2 \subseteq C_1$ , or  $C_1 \cap C_2 = \emptyset$ . The trivial clusters, those contained in every tree with leaf set  $X$ , are  $X$  (the leaf set itself) and the singleton sets  $\{x\}$  where  $x \in X$ . The vertex  $v$  in  $T_3$  of Figure 2.1, for example, defines the cluster  $\{c, d, e, f\}$ , and the set of non-trivial clusters defined by the vertices of  $T_3$  is  $\{\{a, b\}, \{e, f\}, \{d, e, f\}, \{c, d, e, f\}\}$ . We say that the rooted tree  $T'$  refines  $T$  if  $H(T) \subset H(T')$ ; this is denoted  $T \leq T'$ .

The analogous relationships in unrooted trees are bipartitions of the taxon set, and are defined by edges rather than vertices. Each edge  $e$  in a phylogenetic tree  $T$ , when deleted from  $T$ , creates two new subtrees  $T_1$  and  $T_2$ . We say that the edge  $e$  induces the bipartition  $L(T_1)|L(T_2)$ , equivalently



$L(T_2)|L(T_1)$ , of  $L(T)$ ; this bipartition is denoted  $\sigma_e$ . The full set of induced bipartitions of a tree  $T$  is denoted  $\Sigma(T)$ , i.e.  $\Sigma(T) = \{\sigma_e : e \in E(T)\}$ .

The trivial bipartitions, those contained in every unrooted tree with leaf set  $X$ , are those involving a single taxon on one side of the bipartition  $\{x\}|X - \{x\}$  where  $x \in X$ . We often abuse notation, dropping the brackets and commas when writing the sets on either side of the split. Thus, we say that in Figure 2.1, the edge  $(u, v)$  in  $T_1$  induces the bipartition  $ab|cdef$ , and the set of non-trivial bipartitions of  $T_1$  is  $\{ab|cdef, abc|def, abcd|ef\}$ . Naturally, we say that an unrooted tree  $T'$  refines  $T$  if  $\Sigma(T) \subseteq \Sigma(T')$ ; again, this is denoted  $T \leq T'$ . Note that the definition of  $\Sigma(T)$  does not require that  $T$  is unrooted, thus  $\Sigma(T)$  is still defined for rooted trees.  $H(T)$ , on the other hand, does not automatically extend to unrooted trees.

Bipartitions can be used to construct a single tree that summarizes the relationships in a collection of trees with identical leaf sets. The methods used to construct such summary trees are called consensus methods, and the trees produced consensus trees; the most commonly used consensus methods are the strict consensus and majority consensus.

The strict consensus tree of a set of trees is the tree whose set of bipartitions are those that are present in all the trees in the given set. For a given set of trees  $\{T_1, T_2, \dots, T_k\}$ , the strict consensus tree [15]  $T$  is such that

$$\Sigma(T) = \{\sigma : \sigma \in \Sigma(T_i) \text{ for all } 1 \leq i \leq k\} = \bigcap_{i=1}^k \Sigma(T_i) \quad (2.1)$$

The majority consensus tree of a set of trees is the tree whose bipartition set contains those bipartitions present in more than half of the trees in the given set. For a given set of trees  $\{T_1, T_2, \dots, T_k\}$ , the majority consensus tree [16,17]  $T$  is such that

$$\Sigma(T) = \{\sigma : \text{the number of trees } T_i \text{ such that } \sigma \in \Sigma(T_i) > \frac{k}{2}\} \quad (2.2)$$

## 2.1.4 Triplet and Quartet trees

For every three leaves,  $a, b, c$ , there are four possible rooted trees with leaf set  $\{a, b, c\}$ , see Figure 2.3. A rooted binary tree on three leaves is called a triplet. A triplet with leaf set  $\{a, b, c\}$  that contains the cluster  $\{a, b\}$  is denoted  $ab|c$  (equivalently,  $ba|c$ ). The set of triplets displayed by a rooted tree  $T$  is denoted  $r(T)$ ; formally,  $r(T) = \{ab|c : a, b, c \in L(T) \text{ and } \{a, b\} \in H(T|_{\{a,b,c\}})\}$ .

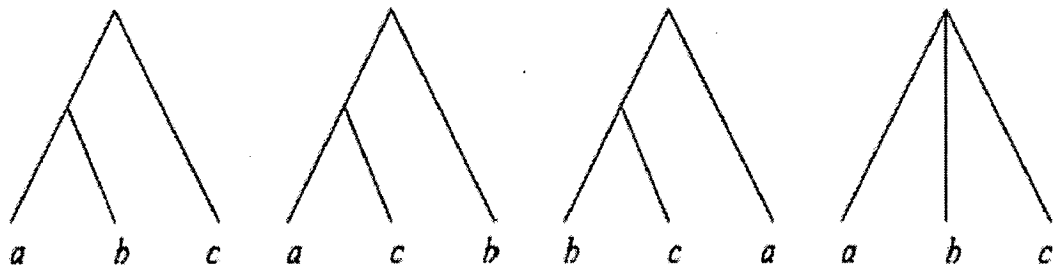


Figure 2.3: The four rooted trees with three leaves.

For every four leaves,  $a, b, c, d$ , there are four possible unrooted trees with leaf set  $\{a, b, c, d\}$ , Figure 2.4. An unrooted binary tree with four leaves is called a quartet tree.

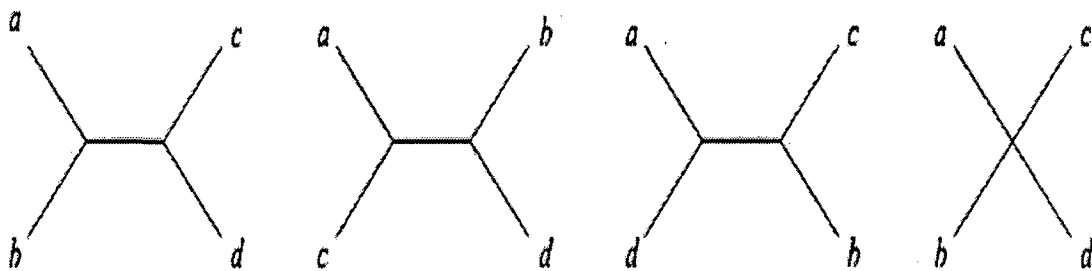


Figure 2.4: The four unrooted trees with four leaves.

A quartet tree with leaf set  $\{a, b, c, d\}$  and non-trivial bipartition  $ab|cd$  is denoted by that bipartition (equivalently,  $ba|cd, ab|dc, ba|dc, orcd|ab$ , etc.). The set of quartet trees displayed by an unrooted tree  $T$  is denoted  $q(T)$ ; formally,  $q(T) = \{ab|cd : a, b, c, d \in L(T) \text{ and } ab|cd = T|_{\{a,b,c,d\}}\}$ .

## 2.2 Phylogenetic Trees

In studying evolutionary histories, we will refer to particular entities of interest, be they species, genes, molecular sequences, or languages, as taxa (singular taxon). This is a generalization of the use of the term taxa in biology where a taxon is any of the taxonomic categories such as phylum, order, or species in the Linnean system.

In general, the taxa have evolved via a process that can be represented by a rooted tree, a directed acyclic graph, in the graph theoretic sense whose root represents the most recent common ancestor of all the taxa represented in the tree. At the leaves of the tree are extant taxa (or possibly extinct taxa for which we have fossil records), the internal vertices are ancestral taxa, and each edge represents an ancestor-descendant relationship with the ancestor being the taxon represented by the vertex closer to the root of the tree.

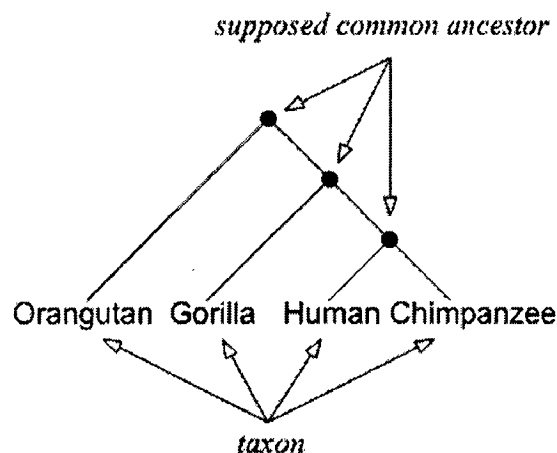


Figure 2.5: Rooted Phylogenetic Tree

Naturally there are times when the assumptions do not hold or are not used, but those are the exceptions. For instance, it might be possible to extract DNA from a specimen that has been preserved somehow, so that some of the taxa might be extinct. Let us take a rooted tree, as like one in Figure 2.5, as an example, systematically giving meaning to the various parts of it. Hopefully it will be clear why rooted trees are ideal for displaying phylogenies, instead of unrooted trees.

- **Root:** The root may be thought of as the (hypothetical) common ancestor of all the taxa.

- **Edges:**An edge corresponds to an ancestor/descendant relationship between the connected nodes: If a node d is a descendant of a, then d has evolved from a somewhere along the line. The length of the edge can be made proportional to the time it has taken to evolve the child from the parent, or some other measure of evolutionary distance. Alternatively, if time is not important, all the edges can have the same length. Edges are directed from ancestor to descendant in a rooted tree, so they can be thought of as time lines.
- **Leaf Nodes:**Leaf nodes represent the existing taxa from which data is available.
- **Internal Nodes:**Each internal node i corresponds to a supposed evolutionary common ancestor of the taxa (leaf nodes) in the subtree rooted by i. Another way to look at it is that i represents a point where a population split up, creating new taxa, which in the tree are represented by the children of i. In any case, the internal nodes are almost always inferred from the biological data of the existing taxa, so we have no real guarantee that such common ancestors ever actually existed in the past.

### 2.3 Phylogenetic supertrees

When researchers publish phylogenetic trees, they are usually the result of a narrowly focused study of some limited taxonomic group. For instance, in several studies of the mammalian order Carnivora, the published phylogenies were rarely larger than 30 taxa. It is simply not viable for a single research team to create large trees. This creates a demand for supertree methods which make it possible to combine the efforts of different teams to create larger phylogenetic trees.

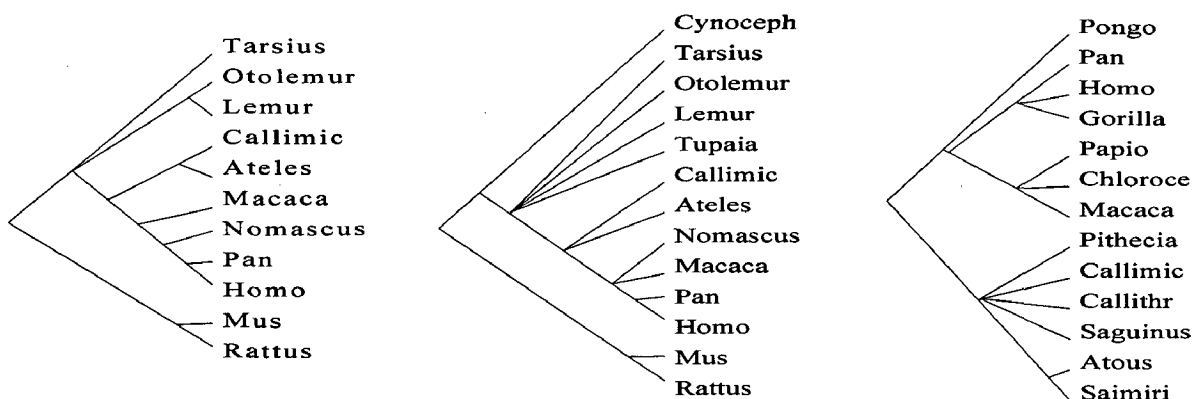


Figure 2.6: An example of input source trees.

There is no need for the teams to use the same methodologies when creating their trees, just as there is no need for them to work on them at the same time. In fact, phylogenetic trees that were published a long time ago, before current supertree methods were invented, can still be used as if they had been published for that very purpose today. Supertrees independence of the method used to create any given phylogenetic tree is a great strength.

Supertrees are themselves phylogenetic trees, but they are built by combining a set of smaller phylogenetic trees. Such a set is referred to as a forest (see Figure 2.6), and the trees in the forest are called source trees.

Regular phylogenetic trees are based on biological data, and this results in the need for a distinction between primary data, which comes directly from the taxa, and secondary data, which is derived from primary data. The source trees must have partially overlapping taxa, or it makes no sense to combine them in a non-trivial fashion. If the source trees are the set  $T = \{T_1, \dots, T_k\}$  then the leaf set of the supertree  $S$  is union of all leaves.

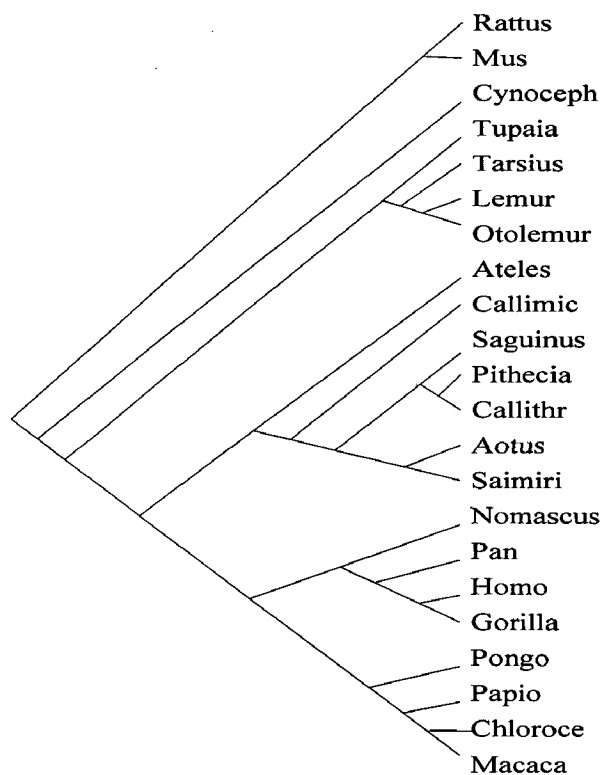


Figure 2.7: A possible supertree for the set of input trees in Figure 2.6.

## **2.4 The Process of a Phylogenetic Study**

When conducting a molecular phylogenetic analysis, biologist's first determine the scope of their study, choosing which phylogenetic group they will study, which particular taxa they will include in the study, and what genetic marker or markers they will sequence to conduct their study. (Some studies may, instead, use morphological characters of the taxa, such as physical characteristics or biological functions, but we will focus on studies that use molecular data.) The limitations that govern which taxa are included in the study include availability of tissue, location of specimens, time, funding, etc., and the ability to amplify and sequence.

### **Phase 1: Data collection**

Biologists go out into the field and collect specimens (organisms or tissues) of the taxa they have decided to study and make use of museum and herbarium specimens. The goal of the data collection phase is to assemble, for each taxon, a homologous sequence or sequences of characters that will be used to construct a phylogenetic tree on the set of taxa they have chosen to study. For studies using molecular characters, these data are molecular sequences for the particular marker they have decided to use in the study. To obtain these sequences, they extract DNA (or RNA) from their specimens and sequence the desired genetic marker. Following sequence assembly, the remainder of the phylogenetic analysis is primarily computational.

### **Phase 2: Sequence Alignment**

An alignment is constructed taking the putatively homologous sequences and putting them in a matrix such that the sequence from each individual taxon occupies its own row, and each column of nucleotides have ideally evolved from a common ancestral nucleotide. This object property of the columns is called positional homology. Positional homology is estimated by inserting gaps into each sequence in the appropriate places, and padding the beginning and end of each sequence with placeholders that represent unknown nucleotides, so as to obtain sequences having the same length. For most phylogenetic reconstruction methods, positional homology is of critical importance because each column of the alignment is assumed to represent the independent evolutionary history of the site. An incorrect alignment will have two or more sites that do not reflect positional homology. Therefore, at least some of the data used to reconstruct the phylogeny will suggest a false evolutionary history of the sequences. To obtain an alignment,

most researchers use some software, such as ClustalW [18], but then often adjust that alignment by eye to create the final alignment.

### Phase 3: Tree Inference

A phylogenetic reconstruction method, such as maximum parsimony, maximum likelihood, Bayesian analysis, or a distance-based method [19], is applied to the aligned sequences. Here we briefly describe maximum parsimony and maximum likelihood:

Maximum Parsimony is an optimization problem based on the minimum evolution principle. A maximum parsimony tree is one that minimizes the number of changes over the edges of the tree needed to “explain” these sequences at its leaves. We formalize this goal here by first defining the parsimony score of a tree.

Given a tree  $T$  that is leaf-labeled by a set of sequences  $S$  each of length  $l$ , an extension  $f$  of  $S$  on  $T$  is a labeling of all vertices of  $T$  by sequences of length  $l$  that maintains the original leaf-labeling by  $S$ . For two sequences  $x$  and  $y$ , let  $H(x, y)$  be the Hamming distance between  $x$  and  $y$ . (For two sequences  $x = x_1x_2\dots x_l$  and  $y = y_1y_2\dots y_l$ ,  $H(x, y) = |\{i : x_i \neq y_i\}|$ .) Then the parsimony score of an extension  $f$ , denoted  $\text{score}(f, T)$ , is  $\sum_{(u,v) \in E(T)} H(f(u), f(v))$ . A minimal extension is one that minimizes the parsimony score, and the parsimony score of a minimal extension is the parsimony length of the tree  $T$ . Computing the parsimony length of a fixed leaf-labeled tree is  $O(nrl)$ , where  $n$  is the number of sequences,  $l$  is the length of the sequences in  $S$ , and  $r$  is the number of states observed in  $S$ . In the case of DNA,  $r = 4$ .

A maximum parsimony tree for a given set of sequences  $S$  is a tree that has the smallest possible parsimony length of any tree leaf labeled by  $S$ . While computing the score of a single tree can be done in linear time, finding a minimum parsimony tree is equivalent to the Hamming-distance Steiner tree problem which is known to be NP-hard [20]. An exhaustive search for a maximum parsimony tree on  $n$  sequences would require looking at all  $(2n - 5)(2n - 7) \dots 1$  possible topologies of trees with the given leaf-labeling. Branch-and-bound algorithms can reduce the search space, but do not reduce the asymptotic running time. Using branch-and-bound algorithms researchers can compute all maximum parsimony trees for datasets of up to about 25 taxa before the running time becomes prohibitive. Thus, for larger datasets, heuristic searches are performed whereby the tree-space is explored until a local optimum is found.

Maximum Likelihood is also an optimization problem, and seeks the most probable tree according to some specified model of evolution. Given a model of sequence evolution, a maximum likelihood tree is a tree topology, along with a set of branch lengths and parameters for the given model of evolution, that maximize the conditional probability of observing the sequences at the leaves given the assumed model. Again let  $S$  be a set of sequences of equal length. The maximum likelihood solution for the set  $S$  is an edge-weighted tree  $(T, w)$ , together with a set of model parameters  $M$ , that minimizes  $\Pr(S|T, w, M)$ . As with MP, ML is NP-hard [21], and heuristic searches are used to compute ML trees.

#### **Phase 4: Post-processing**

In the post-processing stage, scientists summarize the results obtained in the tree inference stage. This typically involves computing a consensus tree for tree inference methods, such as MP, that return multiple trees, and assessing the support for branches in the inferred tree(s).

## **2.5 Literature Review**

UPGMA is a sequential clustering algorithm, and perhaps the simplest method of tree construction. It constructs a rooted tree from a distance matrix  $M$ . First it creates a pool of clusters, each containing a single taxon. Then it combines the two closest clusters in a new cluster  $C$ .  $M$  is updated so that the distance from  $C$  to another cluster  $X$  is the arithmetic mean of the distances from the taxa in  $C$  to  $X$ . This is repeated until there is only one cluster left, at which point the tree can be constructed: Each cluster  $C$  corresponds to a node  $n$ , and the two clusters which comprise  $C$  correspond to the children of  $n$ . The algorithm is criticized for only constructing trees correctly if certain properties are true for  $M$ , namely that  $M$  is ultra-metric.

Neighbour Joining method [22] addresses the main shortcoming of UPGMA by relaxing the requirements for the distance matrix, but it is a linear factor slower and produces unrooted trees. It starts by placing all taxa in a star tree. A modified distance matrix is maintained in which the distance between each pair of nodes is adjusted according to the average distance from all other nodes. In each iteration, the two least distant nodes are linked and removed from the tree, while their common ancestor is added to the tree. The tree is reduced in this way until only two nodes remain, connected by an edge. All the taxa are at this point linked under those nodes, and the tree has been built.



Maximum Parsimony works on the principle that simpler solutions are preferred over more complex ones. For a phylogeny, this means that a tree which can explain the input character data with fewer evolutionary events is preferable over a tree which contains more events. By minimizing the number of character state changes along all paths in the tree, the most parsimonious tree is obtained. Different matrix encoding schemes exist. The original, due to [23, 24], uses 1 to indicate that the trait is present and 0 to indicate that it is not. One of the disadvantages of the method, apart from its NP-hard time complexity, is that no branch lengths are computed: Only the order in which the evolutionary splits occur.

Maximum Likelihood takes as input a model of sequence evolution, a tree providing a topology and branch lengths, and a data matrix. It then computes the likelihood of the tree giving rise to the data matrix, given the specified model of sequence evolution. By searching through the tree space, the trees with the maximum likelihood of being correct can be found. The result depends on the model, but the method is statistically well founded and robust to violations of the rules present in the evolutionary model.

The above methods are used to create phylogenetic trees from primary data, which comes directly from the taxa whereas the following section is about supertree methods which will base on secondary data discussed in section 2.3. Supertree methods summarize the collection of trees without losing branch information. There are different approaches to solve this problem. Some methods solve it using graph theoretic, others are subtree based and yet others are recoding based approach.

Thus supertree methods need to solve the compatibility problem that is to decide if the input trees are compatible. While the compatibility problem can be solved for rooted input trees in polynomial time using the Build algorithm from Aho et al. [25], the problem is NP-complete for unrooted input trees. Thus, the time complexity of supertree problems for rooted and unrooted trees are clearly different. An even stronger separation between supertree methods for rooted and unrooted input trees was made by Steel et al. [26]. Steel et al. showed that in contrast to supertree methods for rooted input trees, no supertree method for unrooted input trees exists that satisfies a subset of the fundamental supertree properties. From this result Semple and Steel [26] conclude that there exists no 'reasonable' supertree method for unrooted input trees.

For any compatible set of triplets and the set of leaf nodes Build algorithm [25] can be used to find the local consensus tree. A fundamental question in biological classification is how a supertree method should use simple input trees to optimally represent their branching information. Probably the oldest supertree algorithm for rooted input trees is the Build algorithm, introduced by Aho et al. [25], and that can be used to solve the compatibility problem in polynomial time. While other supertree problems are typically NP-complete, there are several modifications of the Build algorithm that still run in polynomial time. Semple and Steel [26] introduced a modified version of the Build algorithm, the MinCut Supertree algorithm.

The Build algorithm constructs a rooted tree that is consistent with a given set of lineage constraints. If no such tree exists the empty tree is returned. Given  $n$  lineage constraints the Build algorithm runs in  $O(n^2)$  time [25]. To construct a parent tree for a given profile, the Build algorithm is executed on the set of all lineage constraints that are consistent with a tree in the profile.

Mincut Algorithm is a subtree based algorithm and has all the advantages of using triplets and quartets. Mincut [26] method is one of the very few methods that satisfy all the desirable properties of supertree. This method is an extension of the Build. Build method reports an error when the input trees are incompatible. Semple modified the Build algorithm in such a way that if the input trees are incompatible then some of the edges are removed from the trees based on the minimum loss criteria. The minimum cost is identified using the edges which have minimum weights and when removed will result in a disconnected modified graphs. Therefore the algorithm results in a single tree even if the input trees are not compatible.

The mincut algorithm of [26] was the first to cope with this problem of inconsistent input triplets. Their algorithm handles rooted subtrees of arbitrary size, however, it also specializes in rooted triplets. Page [27] extended this idea to maintain all subtrees that are not in disagreement with any input subtree. The algorithm, denoted modified min cut, applies the min cut criterion but on a somewhat different graph than that of [26].

Modified Mincut method is also a subtree based method. Mincut method results in highly unresolved trees when applied to the input tree representing polytomy. Page [27] modified the Mincut algorithm to work for a general case. The basic idea was to divide the edges of the

modified graph into three categories: unanimous, uncontradicted and contradicted. Whenever the modified graph is fully connected, instead of removing the edges from the minimal cut set of the graph, the contradicted edges are removed. As it is an extension of the Mincut method, it also satisfies the properties that Mincut satisfies.

Indirect supertree construction uses some form of matrix representation (Ponstein 1966, Ragan 1992) to encode individual source tree topologies as matrices that are then combined and analyzed using an optimization criterion. This is the most widely used phylogenetic supertree method defined by Baum [23] and Ragan [24] independently. MRP represented a universally applicable method that could combine the even incompatible set of input trees. In this method binary coding of the components of each input tree is used to generate a pseudo character matrix representation of the trees. The character matrices of all the trees are combined, the leaves that are not present in a given tree are marked as '?', i.e., missing. The final matrix is then analyzed with parsimony procedure [28] to produce one or more parsimonious trees.

Supertree methods developed by Chen et al. [29] and Eulenstein [30] are motivated by the notion of error correction. In MRP taxa present in the same cluster are scored as 1, those absent are scored as 0, and those which are not sampled are marked '?'. One notion of error in such a cluster system is the presence of an incorrect label in a cluster or the absence of one that should be present. This type of errors is called flips  $0 \rightarrow 1$  or  $1 \rightarrow 0$ . This leads to the optimization problem of finding the minimum number of flips that converts the matrix into a matrix which is consistent with a phylogenetic tree. This is an NP-hard problem and [29, 30] gave approximation algorithms for it. Another approach to the matrix representation is to employ the distances between the taxa of each tree and finally combine them into a single average matrix, called as average consensus method [29]. Any distance based tree construction method can then be used to construct the supertree.

## Chapter 3

# Phylogenetic Supertree Construction using Triplet (Inference and) Local Inconsistency

---

### 3.1 Triplet Terminology

For every three leaves let us name,  $a, b, c$ , there are four possible rooted phylogenetic trees with leaf set  $\{a, b, c\}$ , see Figure 3.1. A rooted binary tree on three leaves is called a rooted triplet. A rooted triplet with leaf set  $\{a, b, c\}$  that contains the cluster  $\{a, b\}$  is denoted  $ab|c$  (equivalently,  $ba|c$ ). The set of rooted triples displayed by a rooted tree  $T$  is denoted  $r(T)$ .

Formally,  $r(T) = \{ab|c : a, b, c \in L(T) \text{ and } \{a, b\} \in H(T|_{\{a,b,c\}})\}$ .

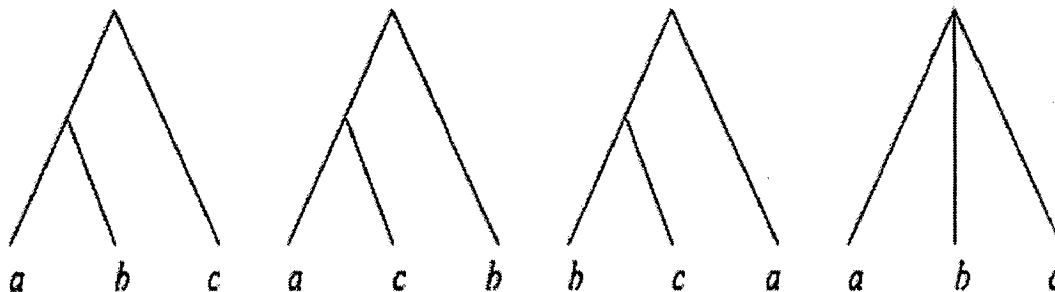


Figure 3.1: Triplet with four non trivial topologies.

In Figure 3.1 for the triplet  $\{a, b, c\}$ , and such a tree is called a triplet tree. A triplet tree is defined uniquely by its ancestral relationships, giving the notation  $a|bc$ ,  $b|ac$ , and  $c|ab$  for the triplet trees displayed in Figure 3.1 from left to right (we ignore the non-binary trifork tree which is unresolved). The lowest common ancestor of two leaf nodes  $a$  and  $b$  is the node  $c$  for which  $c$  is an ancestor of  $a$  and  $b$ , and where no other node with that property is a descendant of  $c$ . We write  $c = lca(a, b)$ , and for a triplet tree  $a|bc$  we have that  $lca(b, c)$  is a descendant of  $lca(a, b)$  and  $lca(a, c)$ .

Triplet tree space is defined as

$$(L) = \bigcup_{M \subseteq L, |M|=3} \{t : t \text{ is a triplet tree on } M\} \text{ on a leaf set } L.$$

The weighting function  $w: T^*(L) \rightarrow R^+$ , and we refer to the triplet trees on a triplet  $M$  as  $t_1^M$ ,  $t_2^M$ , and  $t_3^M$  such that  $w(t_1^M) \leq w(t_2^M) \leq w(t_3^M)$ .

### 3.2 Triplet (Inference and )Local Inconsistency Algorithm [ 31]

The algorithm generates a supertree from input forest of trees by using three steps. In the first step it will decompose the given forest input into set of triplets, second step is an optional one which will be used to infer the missed triplets from the existing triplets and, third will generate a supertree by adding one leaf at a time to the final tree.

```

Algorithm : Decompose [ 31]
Input : forest with leaf set L
Output: a weighting function  $w : T^*(L) \rightarrow R^+$ 
Foreach tree in forest do
    root  $\leftarrow$  root(tree);
    DecorateDown(root);
    DecorateUp(root);
    CountTriplets(root);
end
foreach triplet tree  $t \in T^*(L)$  do
if frequency(t) > 0 then
set  $w(t) = -\log$  frequency(t);
else
set  $w(t) = \infty$ ;
end

```

*Figure 3.2 Algorithm for Decomposing Forest into Triplets*

As said earlier the tree isn't actually reduced to a set of triplet trees, rather it constructs a weighting function  $w: T^*(L) \rightarrow R^+$  for all the triplet trees in the triplet tree space of the leaves of the input forest. We have to count the occurrences of each triplet tree in all the input trees in order to determine their frequencies and thus their weights, which will be used to create the weighting function. The decompose algorithm iterates through the trees of the forest, first decorating them, and then counting the triplet trees in the decorated trees. The algorithms used for each task are all recursive, and their starting point is always the root of the current tree. Once all the trees have had their displayed triplet trees counted, the weighting function is constructed in the above algorithm.

DecorateDown is a recursive algorithm which is the first decoration step and it must be carried out before proceeding, since its results are required for the next decoration step. Each internal node is decorated with a set  $D$  which is to contain the leaf set of the subtree rooted by the node. We shall let  $D_x$  denote the  $D$  set belonging to the node  $x$ . For the input node, first  $D$  is initialized with the empty set. Then all children that are leaf nodes are added to  $D$ , and the  $D$  sets of all the other children, which are internal nodes, are computed recursively and added. Thus,  $D$  contains the leaf nodes directly under the input node, as well as the leaf nodes of the subtrees rooted by its children.

For a node it simply subtracts the leaf nodes in  $D_{\text{node}}$  from the leaf set of the entire tree, and assigns the result to the set  $U_{\text{node}}$ . Then it recurses on those of the children that are internal nodes. The purpose of  $U_{\text{node}}$  is to contain all the leaf nodes which can be reached from the node by traveling at least one edge up, and then an arbitrary path down.

It is possible to create a mapping from the triplets that are resolved in a tree to the tree's internal nodes:

$$\psi: \{M: M \subseteq \text{leaves}(\text{tree}), |M| = 3, \text{tree}|M \text{ resolved}\} \rightarrow \text{internals}(\text{tree}) \quad (4.1)$$

The mapping is defined as follows: For the triplet  $\{a, b, c\}$ , find the lowest common ancestors  $\text{lca}(a, b)$ ,  $\text{lca}(a, c)$ , and  $\text{lca}(b, c)$ , which are all internal nodes. Triplet  $\{a, b, c\}$  is mapped to the node which is the descendant of the others. The node exists because the triplet is resolved otherwise all the lowest common ancestors would be the same.

The mapping tells us that we can count all triplets resolved in a tree by counting the triplets mapped to each internal node, and adding up the results. Since every triplet is mapped to one node, and one node only, no triplet is counted twice, but all triplets are counted.

CountTriplet algorithm receives as input a node. The tree containing node is repeatedly divided into three zones with node in the center. One division corresponds to one pair of node's children (so if the tree is binary, there is only one possible division). The zones correspond to the parts of the tree that contain  $D_x$ ,  $D_y$ , and  $U_{\text{node}}$  for the pair  $x, y \in \text{children}(\text{node})$ .  $U_{\text{node}}$  contains the

candidate nodes for the ancestor leaf nodes in the associated triplet trees, while  $D_x$  and  $D_y$  contain the candidate nodes for the other two leaf nodes.

By picking three leaf nodes, one from each zone, we obtain a triplet. We decide which triplet tree is displayed for the triplet and increment its count. This is repeated for every combination of leaf nodes from the three zones. When all triplet trees have been counted, the algorithm recurses on the children of node that are internal nodes, and thus the whole trees is covered eventually. The time complexity for the above entire procedure is  $O(n^3)$ .

```

Algorithm : Inference [31]
Input : a weighting function  $w : (L) \rightarrow R^+$  for a leaf set L
Output: a modified weighting function  $w : T^*(L) \rightarrow R^+$ 
Data: sets S, R with supported and rejected triplet trees
(S,R)  $\leftarrow$  SortTriplets(w);
Foreach triplet  $M \subseteq L$  do
  Foreach triplet tree  $t \in T^*(M)$  do
    Foreach quartet tree q that displays t do
      For each triplet tree  $p = t$  displayed in q do
        if  $p \in R$  then
          demerit t;
        end
        if  $d \in S$  for some  $d \in T^*(\text{leaves}(p)) \setminus \{p\}$  then
          demerit t;
        end
      end
    end
  end
  if a triplet tree  $t \in T^*(M)$  has fewer demerits than the other then
     $S \leftarrow S \setminus T^*(M)$ ;
     $R \leftarrow R \setminus T^*(M)$ ;
     $S \leftarrow S \cup \{t\}$ ;
  end
end
 $w \leftarrow \text{AdjustWeights}(S,R)$ ;

```

Figure 3.3 Algorithm for inferring missed triplets

SortTriplets function will sort the triplet trees in  $T^*(L)$  into a set of supported and rejected triplet trees called S and R, respectively. Once the triplet trees are sorted, the inference algorithm First, each triplet tree is assigned a demerit counter, which is initially set to zero. Then, the outermost loop iterates through all the triplets  $M \subseteq L$ . In the loop, each triplet tree t for M is examined in

turn to determine its demerits. We can expand  $t$  to a quartet tree by adding a leaf node from  $L \setminus M$ . When expanding a triplet tree, we have four edges where an expansion can take place shown in figure 3.4. A resulting quartet tree displays in 4 Triplets shown in figure 3.5.

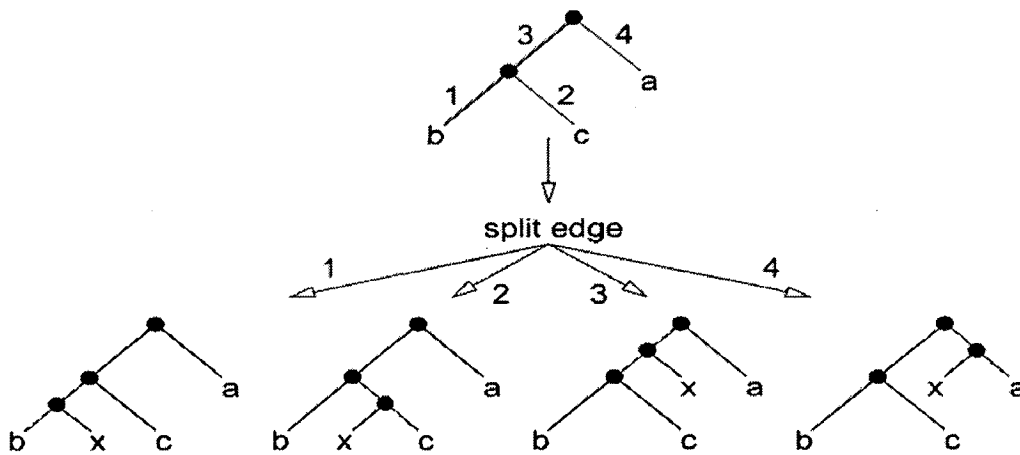


Figure 3.4: An extension of Triplet to a quartet by adding a leaf.

In this case the leaf  $x$  is added in each possible location so there are four possible quartets possible.

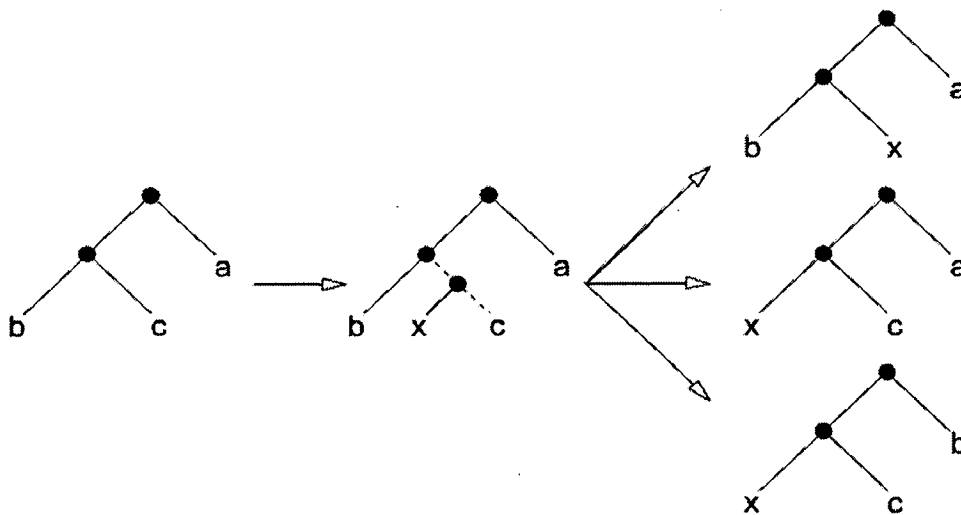


Figure 3.5: A triplet that has been expanded displays three other triplets.



Here  $a|bc$  has been expanded by splitting the edge leading to  $c$ , and the resulting triplet tree displays the additional triplet trees  $a|bx$ ,  $a|xc$ , and  $b|xc$ . The time complexity of the inference step for triplets is  $O(n^4)$ .

```

Algorithm : Compose [31]
Input : a weighting function  $w : T^*(L) \rightarrow R^+$  on the leaf set  $L$ 
Output: a supertree
 $T \leftarrow t_1^M$  such that  $t_1^M$  maximizes  $w(t_2^M) - w(t_1^M)$  over all  $M \subseteq L$ ;
 $L \leftarrow L \setminus \text{leaves}(T)$ ;
while  $L \neq \emptyset$  do
  foreach leaf  $\in L$  do
    foreach tree expanded from  $T$  by adding leaf do
      if  $LI_w(\text{tree}, \text{leaf})$  is the lowest so far then
         $T_2^{\text{leaf}} \leftarrow T_1^{\text{leaf}}$ ;
         $T_1^{\text{leaf}} \leftarrow \text{tree}$ ;
      end
    end
  end
  choose leaf s.t.  $LI_w(T_2^{\text{leaf}}, \text{leaf}) - LI_w(T_1^{\text{leaf}}, \text{leaf})$  is maximized
  and secondarily s.t.  $LI_w(T_1^{\text{leaf}}, \text{leaf})$  is minimized;
   $T \leftarrow T_1^{\text{leaf}}$ ;
   $L \leftarrow L \setminus \{\text{leaf}\}$ ;
end

```

Figure 3.6 Algorithm for composing triplets to a supertree

To determine topology we can use two different approaches, one using properties of lowest common ancestors, and one using a distance matrix. The difference is in their complexities.

### Lowest Common Ancestor (LCA)

Attempt at an algorithm for triplets used the property shown in above, namely that the topology can be detected by comparing lowest common ancestors shown in figure 3.6. First,  $\text{lca}(a, b)$ ,  $\text{lca}(a, c)$ , and  $\text{lca}(b, c)$  are computed. Then we have

$$T|\{a, b, c\} = \begin{cases} a|bc & \text{if } \text{lca}(a, b) = \text{lca}(a, c) \\ b|ac & \text{if } \text{lca}(a, b) = \text{lca}(b, c) \\ c|ab & \text{if } \text{lca}(a, c) = \text{lca}(b, c) \\ abc & \text{lca}(a, b) = \text{lca}(a, c) = \text{lca}(b, c) \end{cases}$$

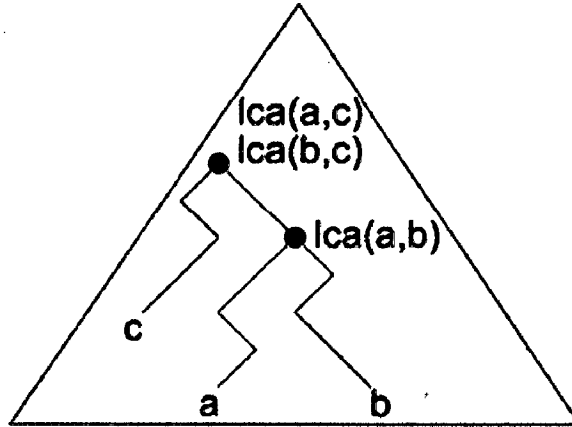


Figure 3.7 the connecting paths between the leaf nodes of a resolved triplet in a tree.

### Distance Matrix (DM)

Distance matrix was based on the idea that the root  $r$  could be added to a triplet as a fourth leaf node. Calculate the path lengths  $D_a|bc = |ra| + |bc|$ ,  $D_b|ac = |rb| + |ac|$ , and  $D_c|ab = |rc| + |ab|$  for each possible topology, and pick the topology which has the shortest one.

we must define excess, inconsistency, and local inconsistency in the following way for a tree  $T$ ; a triplet  $M$ , and a weighting function  $w : T^*(L) \rightarrow R^+$ :

$$excess_w(T, M) = \begin{cases} w(T|M) - w(t_1^M) & \text{if } T|M = t_1^M \\ w(T|M) - w(t_1^M) & \text{otherwise} \end{cases} \quad (4.2)$$

$$I_w(T) = \max \{ excess(T, M) : M \subseteq leaves(T), |M| = 3 \} \quad (4.3)$$

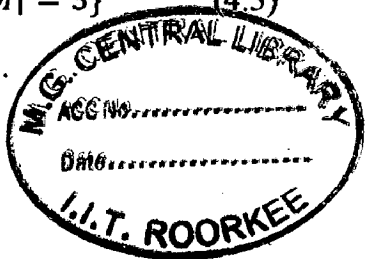
$$LI_w(T, x) = \max \{ excess(T, M) : M \subseteq leaves(T), |M| = 3, x \in M \}.$$

The time complexity for the above algorithm is  $O(n^5)$ .

### 4.3 Implementation and Results

The algorithm Triplet Inference and Inconsistency is implemented in java using MyEclipse IDE. All executions done on a computer with the following specifications

Hardware specifications:



4 GB RAM

Core2Duo processor

Software specifications:

MyEclipse software.

### Results:

The running times are shown in Table 3.1 for TLI using LCA and DM topology detection and TILI using DM topology detection. The executions are on a different forest of the same size. There is an entry for each forest size from 20 to 60 taxa, in increments of ten.

Taxa	TLI(DM)	TLI(LCA)	TILI
20	0.48	0.49	0.78
30	1.79	2.10	3.39
40	6.87	8.50	12.34
50	22.56	29.82	36.62
60	56.27	77.90	85.72

Table 3.1 Time in seconds required to build supertrees from forests of varying sizes.

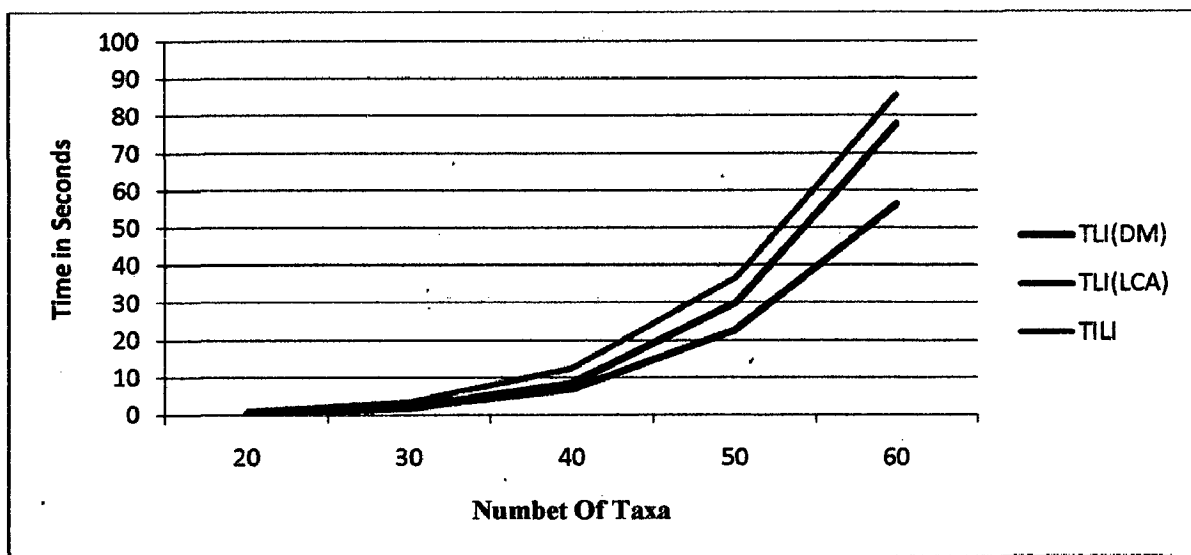


Figure 3.3 Runtime comparison between TLI using both DM and LCA and TILI

From the above we can understand that TLI using DM outperforms the executions times than TLI and TILI. If we use the LCA for computing topology the triplet is taking more time in the excess calculation shown in eq.4.2. Since using LCA for topology detection adds some more time computation for the excess calculation. But the time complexity for computing the topology of a triplet is constant using DM therefore the excess calculation can be done in constant time. For the TILI even using DM it take more time since it involves in computation of the inferring the missing triplets.

## Chapter 4

### Heuristic approach for Phylogenetic Supertree

---

Hill-climbing heuristics have been successfully applied to several intrinsically complex supertree problems [23,32, 33]. They search the space of all possible supertrees guided by a series of exact solutions to instances of a local search problem. The local search problem is to find an optimal phylogenetic tree that attains the optimal criteria in the neighborhood of a given tree. The neighborhood is the set of all phylogenetic trees into which the given tree can be transformed by applying a tree edit operation explained in section 4.2.

In these heuristics a tree graph is defined for the given set of input trees and some, typically symmetric, tree-edit operation. The nodes in the tree graph are the phylogenetic trees over the overall taxon set of the input trees. An edge adjoins two nodes exactly if the corresponding trees can be transformed into each other by the tree edit operation. The cost of a node in the graph is the measurement from the input trees to the tree represented by the node under the particular supertree problems optimization measurement. For the triplet supertree problem defined in later paragraph, the cost of a node in the graph is the triplet-similarity from the input trees to the tree represented by the node. Given a starting node in the tree graph, the heuristic's task is to find a maximal-length path of steepest ascent in the cost of its nodes and to return the last node on such a path. This path is found by solving the local search problem for every node along the path. The local search problem is to find a node with the maximum cost in the neighborhood (all adjacent nodes) of a given node.

Every rooted tree can be equivalently represented by a set of triplet trees [34]. By using this fact we can solve triplet supertree problem i.e. finding a supertree for a given an input profile of  $n$  species trees ( $T_1, \dots, T_n$ ) such that it maximizes the triplet-similarity score which is defined in eq.4.1.

A triplet-similarity measure can be defined between two rooted trees that is the cardinality of the intersection of their triplet presentations. This measure can be extended to measure the similarity from a collection of rooted input trees to a rooted supertree, by summing up the triplet

similarities for each input tree and the supertree. Triplet similarity score between a given profile of trees  $P = (T_1, \dots, T_n)$  and a supertree  $T^*$  of  $P$  can be defined as

$$S(P, T^*) = \sum_{i=1}^n |S(T_i, T^*)| \quad (4.1)$$

The triplet supertree problem is that for a given profile  $P$ , find a supertree  $T^*$  that maximizes  $S(P, T^*)$ . We call any such  $T^*$  a triplet supertree.

## 4.1 Preliminaries

Let  $T$  be a rooted tree. We define  $\leq_T$  to be the partial order on  $V(T)$  where  $x \leq_T y$  if  $y$  is a node on the path between  $\text{Ro}(T)$  and  $x$ . If  $x \leq_T y$  we call  $x$  a descendant of  $y$ , and  $y$  an ancestor of  $x$ . We also define  $x <_T y$  if  $x \leq_T y$  and  $x \neq y$ , in this case we call  $x$  a proper descendant of  $y$ , and  $y$  is a proper ancestor of  $x$ .

The set of minima under  $\leq_T$  is denoted by  $\text{Le}(T)$  and its elements are called leaves. If  $\{x, y\} \in E(T)$  and  $x \leq_T y$  then we call  $y$  the parent of  $x$  denoted by  $\text{Pa}_T(x)$  and we call  $x$  a child of  $y$ . The set of all children of  $y$  is denoted by  $\text{Ch}_T(y)$ . If two nodes in  $T$  have the same parent, they are called siblings. The least common ancestor of a non-empty subset  $L \subseteq V(T)$ , denoted as  $\text{lca}_T(L)$ , is the unique smallest upper bound of  $L$  under  $\leq_T$ .

If  $e \in E(T)$ , we define  $T/e$  to be the tree obtained from  $T$  by identifying the ends of  $e$  and then deleting  $e$ .  $T/e$  is said to be obtained from  $T$  by contracting  $e$ . If  $v$  is a vertex of  $T$  with degree one or two, and  $e$  is an edge incident with  $v$ , the tree  $T/e$  is said to be obtained from  $T$  by suppressing  $v$ . The restricted subtree of  $T$  induced by a non-empty subset  $L \subseteq V(T)$ , denoted as  $T|L$ , is the tree induced by  $L$  where all internal nodes with degree two are suppressed, with the exception of the root node. The subtree of  $T$  rooted at node  $y \in V(T)$ , denoted as  $T_y$ , is the restricted subtree induced by  $\{x \in V(T): x \leq_T y\}$ .

Let  $T$  be a tree and  $v \in V(T)$ , an immediate triplet induced by  $v$ , denoted as  $yz \blacktriangleright \triangleleft v$ , is a triplet  $yz | v$  where there exists nodes  $a, b \in V(T)$  such that  $\text{Pa}_T(y) = \text{Pa}_T(z) = b$  and  $\text{Pa}_T(b) = \text{Pa}_T(v) = a$ .

## 4.2 Tree Edit Operations

A variety of different tree edit operations have been proposed [35], and two of them, rooted Subtree Pruning and Regrafting (SPR) and Tree Bisection and Reconnection (TBR), have shown much promise for phylogenetic studies.

### 4.2.1 ReRoot operation

Let  $T$  be a tree and  $x \in V(T)$ .  $RRT(x)$  is defined to be the tree  $T$  if  $x = Ro(T)$ . Otherwise,  $RRT(x)$  is the tree that is obtained from  $T$  by (i) suppressing  $Ro(T)$ , and (ii) subdividing the edge  $\{Pa_T(x), x\}$  by a new root node.

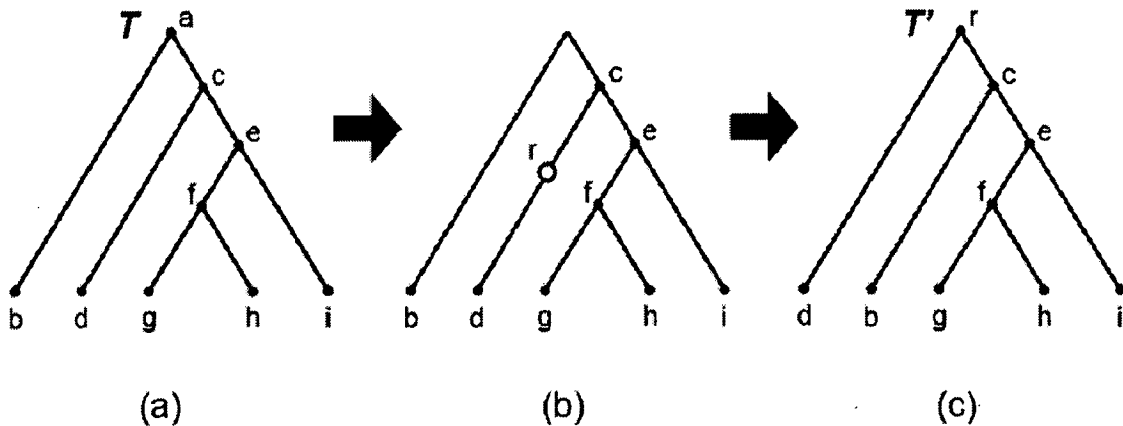


Figure 4.1 Reroot operation

The original tree  $T$  is shown in (a). In (b), we first suppress the root node, and then introduce the new root node  $r$  above  $d$ . Finally we rearrange the tree so that  $r$  is at root, as in (c).

The extension for the RR operation can be defined as  $RR_T = \cup_{x \in V(T)} \{RR_T(x)\}$ .

Let  $x \leq_T v$ , we also define a partial RR operation  $RRT(v, x)$  by replacing  $T_v$  with  $(x)$ .  $RR_{T_v}$

### 4.2.2 Tree Bisection and Reconnection

The planted tree for a tree  $T$  is denoted by  $Pl(T)$  and obtained by adding an additional edge, called root edge,  $\{r, Ro(T)\}$  to  $E(T)$ .

Let  $T$  be a tree,  $e = (u, v) \in E(T)$ , and  $X, Y$  be the connected components that are obtained by removing edge  $e$  from  $T$  where  $v \in X$  and  $u \in Y$ . We define  $TBR_T(v, x, y)$  for  $x \in X$  and  $y \in Y$

to be the tree that is obtained from  $Pl(T)$  by first removing edge  $e$ , then replacing the component  $X$  by  $RR_X(x)$ , and then adjoining a new edge  $f$  between  $x' = Ro(RR_X(x))$  and  $Y$  as follows:

1. Create a new node  $y'$  that subdivides the edge  $(Pa_T(y), y)$ .
2. Adjoin the edge  $f$  between nodes  $x'$  and  $y'$ .
3. Suppress the node  $u$ , and rename  $x'$  as  $v$  and  $y'$  as  $u$ .
4. Contract the root edge.

We say that the tree  $TBR_T(v, x, y)$  is obtained from  $T$  by a tree bisection and reconnection (TBR) operation that bisects the tree  $T$  into the components  $X, Y$  and reconnects them above the nodes  $x, y$ .

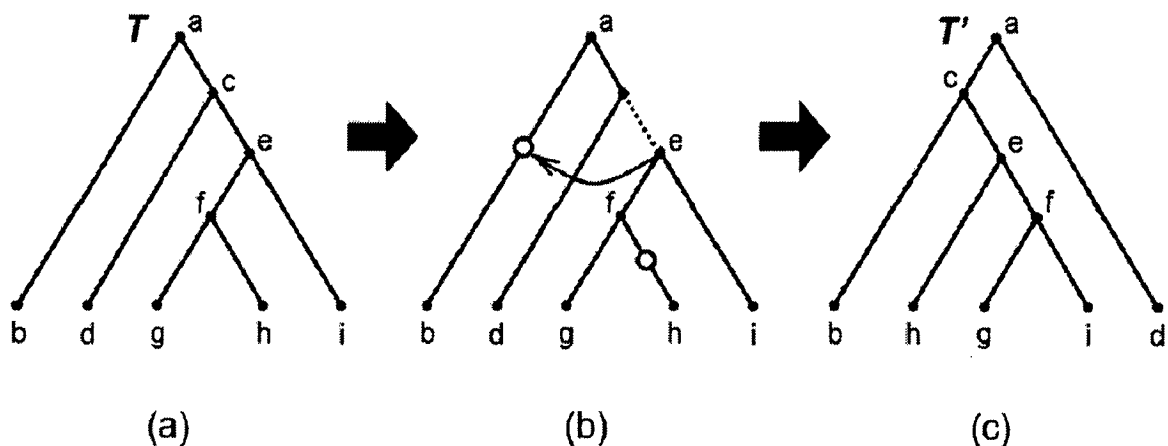


Figure 4.2 Tree Bisection and Reconnect operation.

The original tree  $T$  is shown in (a). In (b), we first remove the edge above  $e$ , that is we prune the subtree  $Te$ . Then we introduce a new node above  $h$  which will be the new root of the pruned subtree. We also introduce a new node above  $b$  this is where we will reconnect the subtree back to  $T$ . Finally we rearrange the tree and obtain the resulting tree  $T'$  as in (c).

We define the following extensions for the TBR operation:

1.  $TBR_T(v, x) = \cup_{y \in Y} TBR_T(v, x, y)$
2.  $TBR_T(v) = \cup_{x \in X} TBR_T(v, x)$
3.  $TBR_T = \cup_{(u, v) \in E(T)} TBR_T(v)$



### 4.1.3 Subtree Pruning and Regrafting operation

An SPR operation for a tree  $T$  can be briefly described through the following steps: (i) prune some subtree  $S$  from  $T$ , (ii) add a root edge to the remaining tree  $T'$ , (iii) regraft  $S$  into an edge of the remaining tree  $T'$ , (iv) contract the root edge.

Let  $T$  be a tree,  $e = (u, v) \in E(T)$ , and  $X, Y$  be the connected components that are obtained by removing edge  $e$  from  $T$  where  $v \in X$  and  $u \in Y$ . We define  $SPR_T(v, y)$  for  $y \in Y$  to be  $TBR_T(v, y)$ . We say that the tree  $SPR_T(v, y)$  is obtained from  $T$  by a subtree prune and regraft (SPR) operation that prunes subtree  $T_v$  and regrafts it above node  $y$ .

We define the following extensions of the SPR operation:

1.  $SPR_T(v) = \cup_{y \in Y} SPR_T(v, y)$
2.  $SPR_T = \cup_{(u, v) \in E(T)} SPR_T(v)$

### 4.2 Triplet Supertree Heuristic Algorithm [36]

SPR Scoring (SPR-S) for a given a profile  $P$  is, a supertree  $T$  of  $P$ , find a tree  $T^* \in SPR_T$  such that

$$S(P, T^*) = \max_{T' \in SPR} S(P, T') \quad (4.2)$$

SPR-Restricted Scoring (SPR-RS) for a given a profile  $P$  is, a supertree  $T$  of  $P$ , and  $(u, v) \in E(T)$ , find a tree  $T^* \in SPR_T(v)$  such that

$$S(P, T^*) = \max_{T' \in SPR_T} S(P, T') \quad (4.3)$$

Algorithm SRR-S problem  $(P, T)$  [36]

Input: A profile  $P = (T_1, \dots, T_n)$ , a supertree  $T$  of  $P$

Output:  $T^* \in SPR_T$ , and  $\Delta_P(T, T^*)$

for all  $(u, v) \in E(T)$  do

Store the value of SPR-RS( $P, T, (u, v)$ )

end for

$(T^*, d) \leftarrow$  the stored value of SPR-RS calls that has the maximum score increase by traversing the tree  $T$  in post order

```

return (T*, d)
end procedure

```

Algorithm for the SPR-RS problem

Input: A profile  $P = (T_1, \dots, T_n)$ , a supertree  $T$  of  $P$ , and  $(u, v) \in E(T)$

Output:  $T^* \in \text{SPR}_T(v)$ , and  $\Delta_P(T, T^*)$

```

r ← Ro(T)

```

```

 $\hat{T} \leftarrow \text{SPR}_T(v, r)$ 

```

```

Call MovedownAndCompute(P,  $\hat{T}$ , v)

```

```

  Traverse the tree  $\hat{T}$  in pre-order to compute  $\Delta_P(T, T')$  for each

```

```

   $T' \in \text{SPRT}(v)$  using the values computed by MovedownAndCompute

```

```

   $T^* \leftarrow T' \in \text{SPR}_T(v)$  such that  $\Delta_P(T, T^*) = \max_{T' \in \text{SPRT}(v)} \Delta_P(\hat{T}, T')$ 

```

```

   $d \leftarrow \Delta_P(\hat{T}, T^*) - \Delta_P(\hat{T}, T)$ 

```

```

  return (T*, d)

```

```

end procedure

```

```

procedure MovedownAndCompute(P, T, v)

```

```

  Input: A profile P, a tree T, and  $v \in V(T)$ 

```

```

   $yz \triangleright \triangleleft v \leftarrow$  The immediate triplet induced by v in T

```

```

  for all  $t \in \{y, z\}$  do

```

```

     $T' \leftarrow \text{SPRT}(v, t)$ 

```

```

    Compute and store  $\Delta_P(T, T')$ 

```

```

    Call MovedownAndCompute(P, T', v)

```

```

  end for

```

```

end procedure

```

*Figure 4.3 Algorithm for SPR-S problem*

Above algorithm is solving the SPR-S problem by using sub algorithm SPR-RS. The SPR-RS algorithm is solving for a given a profile  $P$ , a supertree  $T$  of  $P$ , and  $(u, v) \in E(T)$ , we compute  $\Delta_P(T, T')$  for each  $T' \in \text{SPR}_T(v)$  by first pruning and regrafting  $Tv$  to  $\text{Ro}(T)$  and compute the score differences for each "move-down" operation, then traverse  $T$  in preorder to obtain the tree that gives the maximum score difference.

```

procedure PreprocessTripletSum(P) [ 36]
  Input: A profile P = (T1, ..., Tn)
  Initialize all values of  $\sigma_P$  to 0
  for i = 1 to n do
    for all u  $\in$  V(Ti) in post-order, u  $\notin$  Le(Ti) do
      {v, w}  $\leftarrow$  ChT(u)
      for all {x, y}  $\in$  Le(Tv), z  $\in$  Le(Tw) do
        Increment  $\sigma_P(xy|z)$ 
      end for
      for all {x, y}  $\in$  Le(Tw), z  $\in$  Le(Tv) do
        Increment  $\sigma_P(xy|z)$ 
      end for
    end for
  end for
procedure PreprocessExtendedTripletSum(P, T)
  Input: A profile P = (T1, ..., Tn), a supertree T of P
  for all u  $\in$  V(T) in post-order do
    for all v  $\in$  V(T) in post-order after u, lcaT({u, v})  $\notin$  {u, v} do
      for all w  $\in$  V(T) in post-order after v, lcaT({u, v, w})  $\notin$  {u, v, w} do
        if w  $\in$  Le(T) then
          if v  $\in$  Le(T) then
            if u  $\in$  Le(T) then
               $\sigma_{P, T}(uv|w) \leftarrow \sigma_P(uv|w)$ 
               $\sigma_{P, T}(uw|v) \leftarrow \sigma_P(uw|v)$ 
               $\sigma_{P, T}(vw|u) \leftarrow \sigma_P(vw|u)$ 
            else {u1, u2}  $\leftarrow$  ChT(u)
               $\sigma_{P, T}(uv|w) \leftarrow \sigma_{P, T}(u_1v|w) + \sigma_{P, T}(u_2v|w)$ 
               $\sigma_{P, T}(uw|v) \leftarrow \sigma_{P, T}(u_1w|v) + \sigma_{P, T}(u_2w|v)$ 
               $\sigma_{P, T}(vw|u) \leftarrow \sigma_{P, T}(vw|u_1) + \sigma_{P, T}(vw|u_2)$ 
            end if
          else {v1, v2}  $\leftarrow$  ChT(v)
               $\sigma_{P, T}(uv|w) \leftarrow \sigma_{P, T}(uv_1|w) + \sigma_{P, T}(uv_2|w)$ 
               $\sigma_{P, T}(uw|v) \leftarrow \sigma_{P, T}(uw|v_1) + \sigma_{P, T}(uw|v_2)$ 
               $\sigma_{P, T}(vw|u) \leftarrow \sigma_{P, T}(v_1w|u) + \sigma_{P, T}(v_2w|u)$ 
            end if
          else {w1, w2}  $\leftarrow$  ChT(w)
               $\sigma_{P, T}(uv|w) \leftarrow \sigma_{P, T}(uv|w_1) + \sigma_{P, T}(uv|w_2)$ 
               $\sigma_{P, T}(uw|v) \leftarrow \sigma_{P, T}(uw_1|v) + \sigma_{P, T}(uw_2|v)$ 
               $\sigma_{P, T}(vw|u) \leftarrow \sigma_{P, T}(vw_1|u) + \sigma_{P, T}(vw_2|u)$ 
            end if
          end if
        end for
      end for
    end for
  end for
end for

```

Figure 4.4 Algorithms for Preprocess TripletSum and Extended TripletSum

Above two algorithms are used to preprocess tripletsums. In the preprocess tripletsum the value is assigned to each triplet by the total number of times it present in each tree. For each SPR-S problem Preprocess Extend Tripletsum is called once. The overall time complexity is  $O(n^3)$ . With an expense of initial preprocessing of  $o(kn^3)$ , where  $k$  is number of input trees, and  $n$  is number of taxa in the input trees.

## 4.2 Implementation and Results

The algorithm Triplet Heuristics is implemented in java using MyEclispse IDE. All executions done on a computer with the following specifications

Hardware specifications:

4 GB RAM

Core2Duo processor

Software specifications:

MyEclipse software.

## Results

The running times are shown in Table 4.1 for TILI using DM topology detection and TH . The executions are on a different forest of the same size. There is an entry for each forest size from 20 to 60 taxa, in increments of ten.

Taxa	TILI	TH
20	0.78	0.52
30	3.39	2.18
40	12.34	8.89
50	36.62	20.28
60	85.72	60.25

*Table 4.1 Time in seconds required by TILI and TH for constructing supertree.*

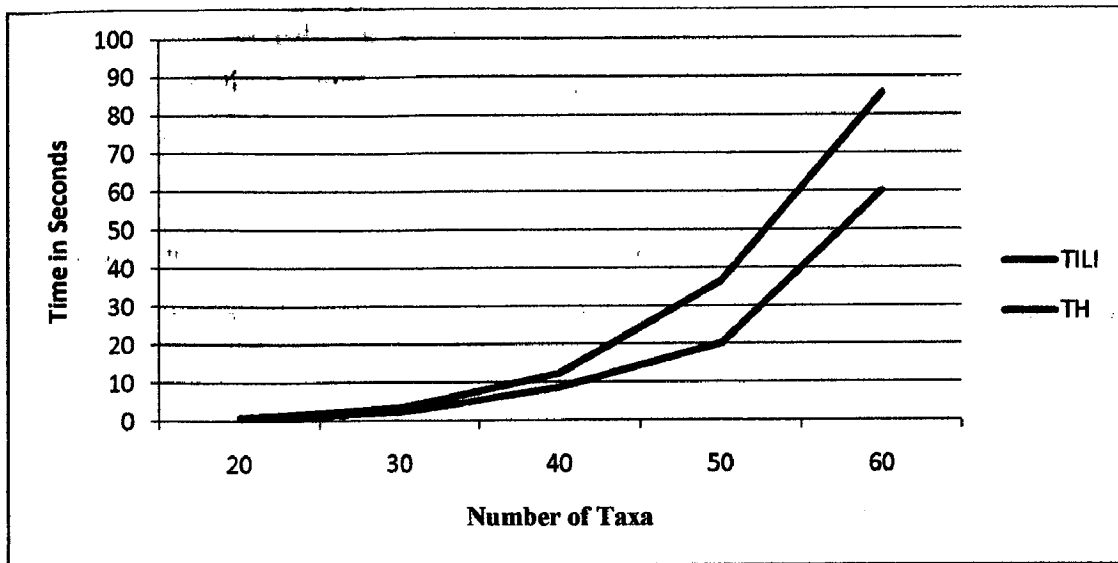


Figure 4.3 Runtime comparisons between TILI and TH

The metric we measure for these two algorithms is run time. From the above we can understand that the running time for construction of supertree for the algorithm TH has less execution time than the algorithm using TILI for the same input set. TILI is involving in inferring missing triplets which results in requiring more execution time.

The other metric we measure is triplet similarity score which is the number of triplets present for the given input profile of trees and the constructed supertree.

Dataset (No .of Taxa)	Method	Triplet similarity
20	TILI	78.12
	TH	97.15
30	TILI	79.24
	TH	97.15
40	TILI	72.15
	TH	98.10
50	TILI	72.15
	TH	98.10
60	TILI	70.17
	TH	98.20

Table 4.2 Triplet Similarity Score between TILI and TH

From the table we can understand that TH algorithm is having more triplet similarity score than TILI algorithm. The triplet similarity for TILI algorithm is reducing as the number of taxa is increasing in the input data set since for each time as the number of taxa increases the input profile of trees is also increases which results in decrease of the triplet similarity for the TILI algorithm. But the TH algorithm is always calculates the similarity score in every step of supertree construction so it is almost constant for it even though the input size increases.

## Chapter 5

### Conclusion and Future Work

---

#### 5.1 Conclusion

The objective of this thesis is implementation of triplet supertree algorithms for construction of phylogenetic supertrees. The following sections describe what has been done:

In this thesis, the requirements for the supertree algorithms are explained. Then the process of phylogenetic study is presented. In chapter 2 the various supertree construction algorithms like Mincut supertree, Modified Mincut algorithm, Matrix Representation with Parsimony are studied. In Matrix Representation with parsimony method, though the solution is guaranteed, the time taken to execute these algorithms will be more compared to both Triplet Inference and Local Inconsistency algorithm (TILI) and Triplet supertree heuristic algorithm (TH).

Later TILI algorithm is implemented in Java +6 SDK and results are obtained under various input data of having different number of taxa. Likewise, TH algorithm is also implemented and executed under same input data. Comparative study based on execution time and accuracy measurement has been done for both the algorithms. The study shows execution time of TH algorithm is less compared to the other one, and is more accurate. From this it can be concluded that TH algorithm is superior to TILI.

#### 5.2 Future Work

Drawing inference is the vital phase in TILI algorithm which in turn affects the accuracy. This can be improved by considering missing triplets through new approaches. This can be further improved by following better weight model.

## REFERENCES

- [1] E. Pennisi. "Modernizing the Tree of Life". *Science*, vol. 300, issue 5626, pp. 1692-1697, 2003.
- [2] T. A. Brown and K. A. Brown. "Ancient DNA: Using molecular biology to explore the past". *Bioassays*, Vol.16, issue 10, pp.719–726,1994.
- [3] B. Kolaczkowski and J. W. Thornton. "Performance of maximum parsimony and likelihood phylogenetics when evolution is heterogeneous". *Nature*, vol.431, pp. 980-984, 2004.
- [4] B. Snel. "Genome trees and the nature of genome evolution". *Annual Review of Microbiology*, vol.59, pp.191-209, 2004.
- [5] A. R. Templeton. "Human origins and analysis of mitochondrial DNA sequences". *Science*, vol.255, issue 5045, pp.737–737, 1992.
- [6] J. Felsenstein. "Phylogenies and the comparative method. *American Naturalist*", vol.125, pp.1–15, 1985.
- [7] D. A. Liberles, D. R. Schreiber, S. Govindarajan, S. G. Chamberlin, and S. A. Benner. "The adaptive evolution database (TAED)". *Genome Biology*, vol.2, issue 8,2001.
- [8] J. A. Eisen. "Phylogenomics: Improving functional predictions for uncharacterized genes by evolutionary analysis". *Genome Research*, vol.8, issue 3, pp.163–167,1998.
- [9] J. R. Brown and P. V. Warren. "Antibiotic discovery: Is it in the genes? *Drug Discovery Today*", vol.3, pp.564–566, 1998.
- [10] R. Overbeek , N. Larsen, and G. D. Pusch. "WIT: integrated system for high throughput genome sequence analysis and metabolic reconstruction". *Nucleic Acids Research*, vol. 28, issue 1, pp.123–125, 2000.
- [11] F. G. Candelas, M. A. Bracho, and A. "Moya. Molecular epidemiology and forensic genetics: Application to a hepatitis c virus transmission event at a hemodialysis unit". *Journal of Infectious Diseases*, vol. 187, issue 3, pp.352–358, 2003.
- [12] J. Gatesy and M. S. Springer, "A critique of matrix representation with parsimony supertrees," in *Phylogenetic Supertrees. Combining Information to Reveal the Tree of Life Computational Biology*, O. R. P. Binnida-Emonds, Ed., vol. 3, pp.369–388, Kluwer Academic Publishers, Dordrecht, The Netherlands, 2004.



- [13] M. Wilkinson, J. A. Cotton, and J. L. Thorley. “*The information content of trees and their matrix representations*”. *Systematic Biology*, vol.53, issue 6, pp.989–1001, 2004.
- [14] O. R. P. Bininda-Emonds. “*MRP supertree construction in the consensus setting*”. In *Bioconsensus*, DIMACS Series in Discrete Mathematics and Theoretical Computer Science, vol.61, pp. 231–242. American Mathematical Society-DIMACS, Providence, Rhode Island, 2003.
- [15] F. R. McMorris, D. B. Meronk, and D. A. Neumann. “*A view of some consensus methods for trees*”. In J. Felsenstein, editor, *Numerical Taxonomy*, pp 122–125. Springer-Verlag, 1983.
- [16] D. L. Swofford. “*When are phylogeny estimates from molecular and morphological data incongruent?*” In M. M. Miyamoto and J. Cracraft, editors, *Phylogenetic analysis of DNA sequences*, pp.295–333. Oxford University Press, 1991.
- [17] J. P. Barthélemy and F. R. McMorris. “*The median procedure for n-trees*”. *J. Classif.*, vol 3, pp.329–334, 1986.
- [18] J. D. Thompson, D. G. Higgins, and T. J. Gibson. “*ClustalW: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice*”. *Nucleic Acids Research*, vol 22, pp.4673–4680, 1994.
- [19] D. L. Swofford, G. J. Olson, P. J. Waddell, and D. M. Hillis. “*Phylogenetic Inference*”, pp 407–425. Sinauer Associates, Sunderland, Massachusetts, 2nd edition, 1996.
- [20] L. R. Foulds and R. L. Graham. “*The steiner problem in phylogeny is NP-complete*”. *Advances in Applied Mathematics*, vol. 3, issue 1, pp.43-49, 1982.
- [21] B. Chor and T. Tuller. “*Maximum likelihood of evolutionary trees: hardness and approximation*”. *Bioinformatics*, vol21, issue 1, pp.97–106, 2005.
- [22] N. Saitou and M. Nei. “*The neighbor-joining method: a new method for reconstructing phylogenetic trees*”. *Molecular Biology and Evolution*, vol 4, pp.406–425, 1987.

- [23] B. R. Baum. “Combining trees as a way of combining data sets for phylogenetic inference, and the desirability of combining gene trees”. *Taxon*, vol.41, pp.3–10, 1992.
- [24] M. A. Ragan. “Matrix representation in reconstructing phylogenetic relationships among the eukaryotes”. *Biosystems*, vol.28, pp.47–55, 1992.
- [25] A. V. Aho, T. G. Sagiv, T. G. Szymanski, and J. D. Ullman. “Inferring a tree from lowest common ancestors with an application to the optimization of relational expressions”. *SIAM Journal on Computing*, vol 10, pp.405–421, 1981.
- [26] C. Semple and M. Steel. “A supertree method for rooted trees”. *Discrete Applied Mathematics*, vol 105, pp.147–148, 2000.
- [27] R. D. M. Page. “Modified mincut supertrees”. *WABI 2002 Lecture Notes in Computer Science*, LNCS-2452, pp.537–551, 2002.
- [28] X. Gu. “Maximum-likelihood approach for gene family evolution under functional divergence”. *Molecular Biology and Evolution*, vol.18, pp.453–464, 2001.
- [29] D. Chen, L. Diao, O. Eulenstein, D. Fernandez-Baca, and M. J. Sanderson. “Flipping: A supertree construction method”. In *Bioconsensus*, pp.135–160. DIMACS series in discrete mathematics and theoretical computerscience, American Mathematical Society, Providence, 2003.
- [30] O. Eulenstein, D. Chen, J. G. Burleigh, D. Fernandez-Baca, and M. J. Sanderson. “Performance of flip-supertree construction with a heuristic algorithm”. *Systematic Biology*, vol .53, issue 2, pp299–308, 2004.
- [31] C. Mosses. “Triplet supertrees”. PhD Thesis, University of Aarhus, Aarhus, Denmark, 2005.
- [32] D. Chen, O. Eulenstein, D. Fernandez-Baca, J. G. Burleigh: “Improved heuristics for minimum-flip supertree construction”. *Evolutionary Bioinformatics*, vol.2, pp.401-410, 2006.

- [33] M. Bansal, O. Eulenstein: "*The Gene-Duplication Problem: NearLinear Time Algorithms for NNI Based Local Searches*". Bioinformatics Research and Applications pp.14-25,2008.
- [34] Bryant D: "*Building Trees, Hunting for Trees, and Comparing Trees – Theory and Methods in Phylogenetic Analysis*". In PhDThesis University of Canterbury; 1997.
- [35] Semple C, Steel M: "*Phylogenetics*", Oxford University Press; 2003.
- [36] H. T. Lin, J. G. Burleigh , O. Eulenstein: "*Triplet supertree heuristics for the tree of life*". BMC Bioinformatics, vol.10, issue 1, pages:S8,2009.