

On a Mirkin-Muchnik-Smith Conjecture for Comparing Molecular Phylogenies*

Louxin Zhang[†]

Abstract

A conjecture of Mirkin, Muchnik and Smith is answered affirmatively which connects the inconsistency function, a biologically meaningful similarity/dissimilarity measure for a gene tree and a species tree, to the mutation cost function, a combinatorial measure based on the mapping of trees. A linear-time algorithm for computing the mutation cost function is also derived from the conjecture.

1 Introduction

As DNA sequences have become easier to obtain, interesting emphasis has been placed on constructing gene trees and from these, reconstructing evolutionary trees for species. Because of the presence of paralogy and sorting of ancestral polymorphism, gene trees and species trees are often inconsistent (e.g., Neigel and Avise, 1986; Pamilo and Nei, 1988; Takahata, 1989; Wu, 1991). Therefore, a major concern that arises is how to combine different, sometimes contradictory, gene trees into an evolutionary tree called species tree (Fitch, 1970; Goodman *et al.*, 1979; Nei, 1987). Several ideas have been suggested for the last twenty years (see, for example, Robinson 1971; Waterman and Smith, 1978; Margush and McMorris, 1981; Hendy *et al.*, 1984; Adams, 1986; Barthélemy *et al.*, 1986). A common characterization of these ideas is that they consider phylogenetic trees as formal mathematical objects and proposed the similarity/dissimilarity measures based only on combinatorial consideration. The weakness of these measures is lack of biological meaning (Mirkin *et al.*, 1995) and computable intractability in general (see, for example, DasGupta *et al.*, 1996).

Biologically meaningful similarity/dissimilarity measures are also addressed. Since gene divergence causes all contradictions among different gene trees, it should be presented and explained in the combined species tree (Fitch, 1970; Goodman *et al.*, 1979). The gene divergence can be the results of either speciation or duplication (Ohno, 1970). The speciation happens between species. If the gene divergence occurs with only speciation, the gene and species trees are identical. But, the duplication event happens within species. When duplications occur, the gene and species trees might be inconsistent. If the common ancestry of two genes

* Appeared in **J. of Comput. Bio. Vol. 4 (1997), No. 2, pp. 177-187**

[†]Supported by a CGAT grant. This work was done when the author was in the Department of Computer Science, University of Waterloo. Current address: BioInformatics Center, Institute of Systems Science, Heng Mui Keng Terrace, Singapore 119597. E-mail: lxzhang@iss.nus.sg

can be tracked back to a speciation event, then they are said to be related by *orthology*; if it is tracked back to a duplication event, then they are related by *paralogy* (Fitch, 1970). Taking into account orthology and paralogy evolutions, Goodman *et al.* (1979) proposed a new similarity/dissimilarity measure for annotating species tree with duplications, gene losses and the nucleotide replacements. Later, Guigó *et al.* (1994) elaborated the idea for identifying and locating the gene duplications in eukaryotic history.

Guigó *et al.* (1994) introduced the mutation cost functions for measuring the dissimilarity between gene and species trees using so called the mapping of trees. The mapping of trees was considered implicitly in Goodman *et al.* (1979) and explicitly in Page (1994). As we shall see later, the concept of mapping and thus the mutation cost function are rather formal and technical. Therefore, the definitions of duplications and subsequent events based on mapping are not substantiated with any biologically meaningful model. This leads Mirkin *et al.* (1995) to propose a new, biological meaningful model for explaining and measuring the dissimilarity between a single gene tree and a species tree. They first formalized the concepts of gene duplication, loss, information gap in graph-theoretic terms, and then proposed a procedure for measuring the dissimilarity between a gene and species tree by comparing every subtree of the gene tree with the corresponding subtree in the species tree. Whenever an inconsistency occurs, a duplication event is assumed to explain the inconsistency. This duplication is reflected in the species tree with the event's history leading to the current situation in species. The history along the tree involves certain gene copy losses which accompany the duplication during the evolutionary history. The total number of duplications, loss and information gap events involved in all the inconsistency is defined to be the inconsistency measure between the gene and species tree. In the same paper, Mirkin *et al.* conjecture that the inconsistency measure coincides with the mutation cost function.

The major goal of our work is to prove that the conjecture is true. This may justify the use of the mutation cost function in evolutionary tree reconstruction. Our work also provides a linear-time algorithm for computing the mutation cost function.

The rest of the paper is divided into two sections. In Section 2, we briefly review the biologically consistency model of Mirkin *et al.*, and its basic properties. In Section 3, we show that the above mentioned conjecture of Mirkin *et al.* is true. Then we present linear-time algorithms for computing mapping from a gene tree to a species tree and the mutation cost function.

2 A model based on gene duplications

In this section we briefly introduce the Mirkin-Muchnik-Smith model for comparing a gene tree and a species tree. For its biological meaning, we refer the reader to Mirkin *et al.* (1995).

2.1 Gene duplications

For a set I of N biological taxa, the model for their evolutionary history is a full, rooted binary tree T with N leaves each labeled by a distinct element of I . Any internal node denotes an ancestor of some taxa in I and are considered as a subset (also called cluster) of its subordinate leaves. Thus, the evolutionary relation “ m is a descendant of n ” is expressed, in set-theoretic

setting, just as “ $m \subset n$ ”, where we use strict inclusion, in contrast to notation $m \subseteq n$, which allows the equality of m and n .

Each internal node has two children, which are denoted by $a(n)$ and $b(n)$. If n_1, n_2, \dots, n_l is a path connecting node n_1 and its descendant n_l , then $n_l \subset n_{l-1} \subset \dots \subset n_1$, and any node m belongs to the path between n_1 and n_l if and only if $n_l \subset m \subset n_1$, which is called an *intermediate* between n_1 and n_l .

A subset S of nodes of T is incompatible if $x \cap y = \phi$ for any $x, y \in S$. For an incompatible subset S in T , we denote by $T|S$ the smallest subtree T containing S as its leaf set. The *homomorphic subtree* of T induced by S is the subgraph obtained from $T|S$ by contracting all degree 2 nodes except for its root. These concepts are illustrated in Figure 1. Finally, we denote the root of T by $r(T)$.

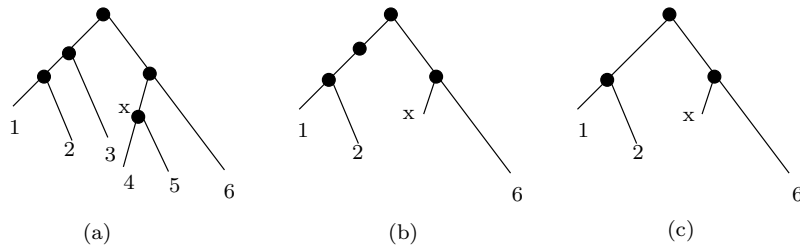


Figure 1: (a) A phylogenetic tree T ; (b) The subtree $T|S$ for $S = \{1, 2, x, 6\}$; (c) The homomorphic subtree induced by S .

In order to annotate duplication history into a phylogenetic tree, the concept of duplication is introduced in Mirkin *et al.* (1995).

Definition 2.1 Let T be a phylogenetic tree with leaf set I representing a set of N species and let $L = \{\phi, +, -, +-\}$. A mapping $\delta : T \rightarrow L$ is called a gene duplication if (1) it is monotone, that is, $\delta(m) \subseteq \delta(n)$ when $m \subseteq n$, and, (2) it is saturated, that is, $\delta(t) = +-$ for the root t of T .

One could think that the duplication δ emerges in the root n and then evolves through all its descendants by only losing sometimes, one (speciation) or both (lack of information) gene copies, which are represented by the signs $+$ and $-$.

A node $m \in T$ is called *mixed* if $\delta(m) = +-$, *speciated* if $\delta(m) = +$ or $-$, and *gapped* if $\delta(m) = \phi$. A node is *speciated/gapped* if it is either speciated or gapped. The maximal mixed, speciated/gapped nodes have particular evolutionary meaning: the maximal mixed node is the root, which represents the duplication event itself, the maximal speciated/gapped nodes correspond to the gene losses. The total number of these maximal nodes will be called the *complexity* of duplication δ , which counts for the total number of the evolutionary distinct events associated with the duplication¹.

¹In Mirkin *et al.* (1995), the complexity is differently defined as the total number of maximal duplications, maximal speciated and maximal gapped nodes. For so-called operational duplications, their definition and ours become identical. Main reason for modifying the definition is that under our definition, the conjecture is always true, not just for operational duplications.

2.2 Comparing a gene tree with a species tree

Let T be a species tree rooted at t with leaf set I and let G be a single gene tree rooted at g with leaf set J such that $J \subseteq I$. T and G will be said to be *root-consistent* if each of the sets of the leaf descendants, $a(g), b(g)$, of the children of the root g is contained in a child-set $a(t)$ or $b(t)$ of t , and *root-inconsistent* otherwise. Obviously, the root-consistency means that the ‘root branches’ of the tree G consistent with the ‘root branches’ of T ; an event causing the divergence of the root descendants occurring in the species tree T is also reflected in the divergence of the specific gene family represented by the gene tree G .

If G and T are root-inconsistent, we use a duplication event in the root of T to express the inconsistency. Thus, we pose the following postulate.

Duplication/Speciation Principle(Mirkin *et al.*, 1995). *Root-inconsistency of the trees G and T means that a duplication event in the gene corresponding to tree G occurs at the root of the species tree T and evolves in T in such a way that contemporary organisms corresponding to the leaves in $a(g)$ have one copy of the divergent gene and the leaves in $b(g)$ the other.*

Duplication/Speciation Principle induces the following simplest duplication assignment $\delta_g : T \rightarrow \{+-, +, -, \phi\}$:

1. for any leaf $i \in I$, $\delta_g(i)$ is defined by

$$\delta_g(i) = \begin{cases} + & \text{if } i \in a(g), \\ - & \text{if } i \in b(g), \\ \phi & \text{if } i \in I - J; \end{cases}$$

2. for any internal node $n \in T$,

$$\delta_g(n) = \begin{cases} + & \text{if } n \cap a(g) \neq \phi \text{ and } n \cap b(g) = \phi, \\ - & \text{if } n \cap b(g) \neq \phi \text{ and } n \cap a(g) = \phi, \\ +- & \text{if } n \cap a(g) \neq \phi \text{ and } n \cap b(g) \neq \phi, \\ \phi & \text{if } n \subseteq I - J. \end{cases}$$

The mapping δ_g can be computed just in one bottom-up run through tree T using at most N steps, each involving comparing two sets of at most N elements. The complexity $c(\delta_g)$ demonstrates the extent of the biologically meaningful difference between G and T , and will be denoted as $c(g, T) = c(r(G), T)$.

To compare the entire gene tree G with species tree T , we need comparing all the subtrees of G with those of T . For simplicity, we assume that these two trees have an identical leaf set, that is, $I = J$. Recall that $T(n)$ denotes the subtree of T rooted at the node n . The comparing procedure consists of sequential comparisons of all the gene subtrees $G(m)$ with those subtrees, $T(n)$, of T such that $m \subseteq n$. If $G(m)$ is root-consistent with a subtree $T(n)$ of T , we proceed to subtrees of $T(n)$. When $G(m)$ is root-inconsistent with $T(n)$, the minimum duplication assignment $\delta_{nm} : T(n) \rightarrow \{+-, +, -, \phi\}$ is defined as a duplication of gene G in the node n ($\in T$) to explain the inconsistency. Figure 2 illustrates the comparison of a gene tree G (in (b)) with a species tree T (in (a)). Three subtrees rooted at A, B , and C are root-inconsistent with the corresponding subtrees G , $\{4, 5, 6\}$ and $\{7, 8, 9\}$ respectively. The

corresponding duplications are shown in Figure 2 (c), (d), and (e) respectively, where maximal speciated/gapped nodes are marked with square boxes. The costs of duplications in (c), (d), and (e) are 8, 5 and 4 respectively.

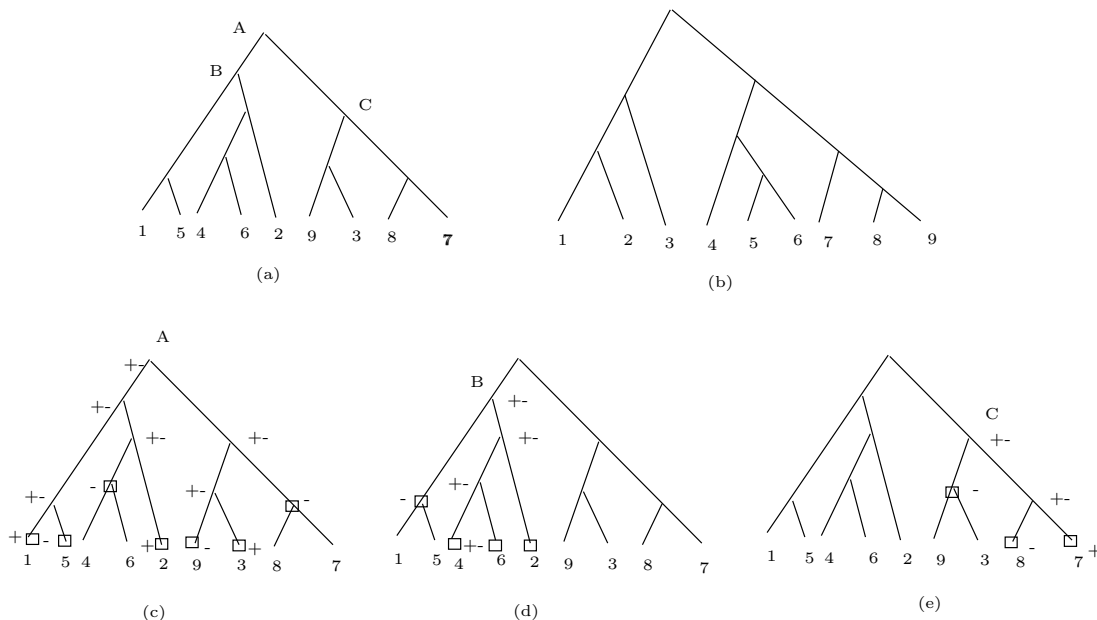


Figure 2: *Duplications between a species tree (a) and a gene tree (b)*

The total inconsistency, $c(T/G)$, between T and G is defined through all the duplication events:

$$c(T/G) = \sum_{g' \in G} c(g', T),$$

where $c(g', T) = 0$ if there is no duplications involving g' in T .

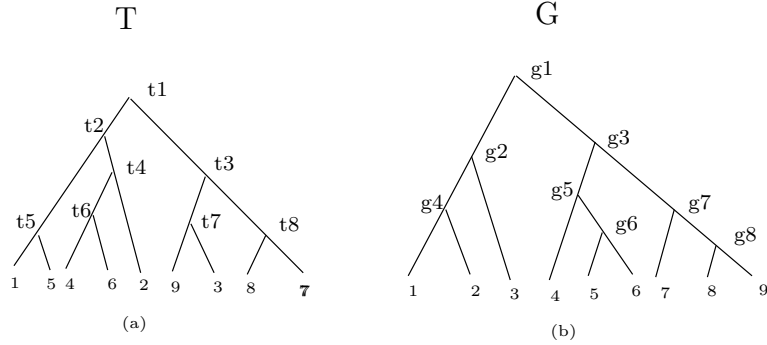
3 Computing the inconsistency

3.1 Mirkin-Muchnik-Smith conjecture

Given a species tree T and a gene tree G with the same leaf set of N taxa. By definition, computing the total inconsistency requires $O(N^3)$ steps: any of $(N - 1)$ non-singleton subtrees of G is compared with at most $(N - 1)$ subtrees of T , and each comparison involves checking root-inconsistency and then defining the duplication if necessary, which takes at most N^2 steps.

Actually, computing the total inconsistency of T and G is much easier based on a conjecture posed by Mirkin *et al.* (1995) which relates $c(T/G)$ to the combinatorial properties of mapping G into T . Before stating the conjecture, we introduce some necessary concepts and facts regarding to mapping G into T .

For any node $g \in G$, we use $M(g)$ to denote the node of T being its least common ancestor, that is, the smallest cluster containing g . This correspondence M , first considered by Goodman



Node and its children	Destinations	Duplication event	Cost
$\begin{pmatrix} g1 \\ g2 \\ g3 \end{pmatrix}$	$\begin{pmatrix} t1 \\ t1 \\ t1 \end{pmatrix}$	✓	1
$\begin{pmatrix} g2 \\ g4 \\ 3 \end{pmatrix}$	$\begin{pmatrix} t1 \\ t2 \\ 3 \end{pmatrix}$		2
$\begin{pmatrix} g4 \\ 1 \\ 2 \end{pmatrix}$	$\begin{pmatrix} t2 \\ 1 \\ 2 \end{pmatrix}$		2
$\begin{pmatrix} g3 \\ g5 \\ g7 \end{pmatrix}$	$\begin{pmatrix} t1 \\ t2 \\ t3 \end{pmatrix}$		0
$\begin{pmatrix} g5 \\ 4 \\ g6 \end{pmatrix}$	$\begin{pmatrix} t2 \\ 4 \\ t2 \end{pmatrix}$	✓	4
$\begin{pmatrix} g6 \\ 5 \\ 6 \end{pmatrix}$	$\begin{pmatrix} t2 \\ 5 \\ 6 \end{pmatrix}$		3
$\begin{pmatrix} g7 \\ 7 \\ g8 \end{pmatrix}$	$\begin{pmatrix} t3 \\ 7 \\ t3 \end{pmatrix}$	✓	3
$\begin{pmatrix} g8 \\ 8 \\ 9 \end{pmatrix}$	$\begin{pmatrix} t3 \\ 8 \\ 9 \end{pmatrix}$		2

Figure 3: Mapping G onto T .

et al. (1979), is referred to as mapping of G into T by Page (1994). We call $M(g)$ the *destination* of g . Obviously, if $g' \subset g$, then $M(g') \subseteq M(g)$, and any leaf is mapped onto a leaf with the same label. Recall that for an internal node g , $a(g)$ and $b(g)$ denotes its two children.

Definition 3.1 Let g be an internal node of G . It is said to be *type-1* under the mapping if $M(a(g)) \subset M(g)$ and $M(b(g)) \subset M(g)$; it is *type-2* if $M(a(g)) \subset M(g)$ and $M(b(g)) = M(g)$ or vice versa; it is *type-3* if $M(a(g)) = M(b(g)) = M(g)$.

Obviously, the following fact is true.

Proposition 3.1 (Mirkin et al., 1995) A node $g \in G$ is a *type-2* or *type-3* node if and only if the subtrees $G(g)$ and $T(M(g))$ are root-inconsistent.

The mutation cost function associated with the mapping itself is reformulated as follows. We use G_i to denote the set of all *type- i* nodes in G for $i = 1, 2, 3$.

Definition 3.2 (Guigó et al., 1994) The cost $L(g)$ associated with $g \in G$ is defined as:

$$L(g) = \begin{cases} |M(g)M(a(g))| + |M(g)M(b(g))| & \text{if } g \in G_1 \\ |M(g)M(a(g))| + 2 & \text{if } g \in G_2 \text{ and } M(a(g)) \subset M(g), \\ 1 & \text{if } g \in G_3. \end{cases}$$

where $|M(g)M(a(g))|$ denotes the number of intermediate nodes between $M(g)$ and its descendant $M(a(g))$ in T .

The mutation cost function $c(G, T)$ associated with the mapping of G into T is the sum of all penalties $L(g)$ over all the internal nodes $g \in G$.

Figure 3 presents the mapping of the gene tree G onto the species tree T which are illustrated in Figure 2. Obviously, when G and T are identical, $c(G, T) = 0$. In their paper, Mirkin *et al.* (1995) conjecture that

$$c(T/G) = c(G, T)$$

for any two trees with the same labeled leaf set.

3.2 The proof of the conjecture

In this section, we shall prove that their conjecture holds. Here, we assume that T is a species tree and G a gene tree and both have the same leaf set I . First, we introduce a simple fact, which appears implicitly in the proof of Theorem 9 in Mirkin *et al.* (1995). However, for completeness, we give its proof.

Proposition 3.2 *If T and G are root-inconsistent and δ denotes the duplication assignment from G to T induced by the Duplication/Speciation Principle, then the number of the maximal speciated/gapped nodes is equal to the number of the mixed nodes plus 1.*

Proof. Let t be the root of T . Recall that T is a full binary tree. Consider a maximal speciated/gapped node $n \in T$ in the assignment δ . Since it is maximal, any intermediate node between n and the root of T must be a mixed node. Conversely, since for each leaf $l \in T$, $\delta(l) = +, -, \phi$, there is a maximal speciation/gapped node on the path from t to l . Let T' denote the restriction of T on all mixed and maximal speciated/gapped nodes. Then T' is a subtree with the same root t . Its leaf set consists of all maximal speciated/gapped nodes in T , while its all internal nodes are mixed under δ . Since T is full, each internal node of T' has exactly two children in T' and so the number of its leaves is equal to 1 plus the number of its internal nodes. \square

Now we consider the duplication mapping M from G to T . Suppose that under the mapping M , there are k_1 type-1 nodes,

$$g_{11}, g_{12}, \dots, g_{1k_1},$$

k_2 type-2 nodes,

$$g_{21}, g_{22}, \dots, g_{2k_2},$$

and k_3 type-3 nodes,

$$g_{31}, g_{32}, \dots, g_{3k_3},$$

in G respectively. Since G is a full binary tree with N leaves, G has $N - 1$ internal nodes. Therefore,

$$k_1 + k_2 + k_3 = N - 1. \quad (1)$$

By Proposition 3.1, we also have

Proposition 3.3 *There are $k_2 + k_3$ duplications between G and T .*

For a type-1 node g_{1i} , $1 \leq i \leq k_1$, $M(a(g_{1i}))$ and $M(b(g_{1i}))$ are distinct from $M(g_{1i})$. The unique path from $M(a(g_{1i}))$ to $M(b(g_{1i}))$ through $M(g_{1i})$ is called an *arc* in the mapping M from G to T . For our purpose, we say that such an arc *starts* at $M(g_{1i})$. We also say that such an arc *passes* through any intermediate between $M(a(g_{1i}))$ and $M(g_{1i})$ or between $M(b(g_{1i}))$ and $M(g_{1i})$. For a type-2 node g_{2i} , $1 \leq i \leq k_2$, let $M(a(g_{2i})) \subset M(g_{2i})$ and let $M(b(g_{2i})) = M(g_{2i})$. The unique path from $M(g_{2i})$ to its descendant $M(a(g_{2i}))$ is called a *path* in the mapping M from G to T , *starting* at $M(g_{2i})$. Such a path *passes* through all intermediates between $M(a(g_{1i}))$ and $M(g_{1i})$.

Proposition 3.4 *For any non-duplication node $x \in T$, the total number of duplications (between G and T) in which x is mixed is exactly one less than the number of arcs and paths passing through x in the mapping M of G in T .*

For any duplication node $x \in T$, the total number of duplications in which x is mixed is one less than the sum of the numbers of the arcs and paths passing through x and of the arcs and paths starting at x .

For understanding Proposition 3.4 and the following proof, one would better study the example illustrated in Figure 2 and Figure 3 again.

Proof. Since the first statement is just a special case of the second, we just prove the second. Consider a node $x \in T$. Set

$$\begin{aligned} G'_x &= \{g' \in G \mid M(g') \subset x \subset M(p(g'))\}, \\ G''_x &= \{g'' \in G \mid x \cap g'' \neq \phi \ \& \ x \subseteq M(g'')\}, \\ \overline{G}'_x &= \{g' \in G \mid M(g') \cap x = \phi \ \& \ a(p(g')) \in G'_x\}, \end{aligned}$$

where $a(p(g'))$ is the other child of $p(g')$. It is not difficult to see that G''_x consists of the root $r(G)$ of G and intermediate nodes between the nodes in G'_x and $r(G)$, and \overline{G}'_x consists of siblings of nodes in G'_x . Moreover, by definition, the parents of all nodes in G'_x are either type-1 or type-2 nodes and the corresponding paths starting at their destinations pass through x . In fact, as we shall see later, for a node $g \in G'_x$, if its sibling is in \overline{G}'_x , then its parent is a type-1 node (Claim 3); if its sibling is in G''_x , then its parent is type-2 (Claim 4). Therefore, $G'_x \cup G''_x$ is a subtree rooted at $r(G)$. Consider the homomorphic subtree, G_x , of G induced by nodes in G'_x . Recall that G_x is obtained from $G'_x \cup G''_x$ after the contraction of all degree-2 nodes except for the root, which is $r(G)$.

Claim 1. Let $y \in G$. If y induces a duplication at $M(y) \in T$ in which x is mixed, then y is a node in the subtree G_x .

Proof. Since duplication occurs in $M(y)$, then $G(y)$ and $T(M(y))$ are root-inconsistent. Therefore, $M(a(y)) = M(y)$ or/and $M(b(y)) = M(y)$. Since x is mixed in this duplication, we

have $x \subseteq M(y)$, which implies that $y \in G_x''$. Since $a(y) \cap x \neq \phi$ and $b(y) \cap x \neq \phi$, y is either the root of $G_x' \cup G_x''$ or a degree-3 node in $G_x' \cup G_x''$, and thus is in G_x . \square

Conversely, we have

Claim 2. Let $y \in G_x$. If $a(y) \in G_x'' \cup G_x'$ and $b(y) \in G_x'' \cup G_x'$, then y is a type-3 or type-2 node and x is mixed in the duplication induced by y .

Proof. The case both $a(y) \in G_x'$ and $b(y) \in G_x'$ is impossible. We consider other three cases as follows. (1) If $a(y) \in G_x''$ and $b(y) \in G_x''$, then $x \subseteq M(a(y))$ and $x \subseteq M(b(y))$, and both $M(a(y))$ and $M(b(y))$ are in the path from the root t to x . If $M(b(y)) \subseteq M(a(y))$, then $M(y) = M(a(y))$, and $M(y) = M(b(y))$ otherwise. Thus, y is a type-2 or type-3 node. By the definition of G_x'' , $x \cap a(y)$ and $x \cap b(y)$ are non-empty, and so x is mixed in the duplication induced by y . (2) If $a(y) \in G_x''$, but $b(y) \in G_x'$, then $M(b(y)) \subset x$ and $x \subseteq M(a(y))$. We have $M(a(y)) = M(y)$, and thus y is a type-2 node. Moreover, $x \cap a(y)$ and $x \cap b(y)$ are non-empty. Therefore, x is mixed in the duplication induced by y . (3) If $b(y) \in G_x''$, but $a(y) \in G_x'$, we can prove th y is also a type-2 node as the second case. \square

Claim 1 and Claim 2 reflect that all duplications in which x is mixed are induced by the corresponding internal nodes of G_x . On the other hand, there are also one-to-one corresponding between all arcs and paths passing through x and all type-1 or type-2 nodes in G_x . Given a node $y \in G$, suppose y is a type-1 or type-2 node and the corresponding arc or path passes through x . By the construction of G_x , $y \in G_x$. Furthermore, we have the following facts.

Claim 3. Let $y \in G_x$. If $a(y) \in G_x'$ and $b(y) \in \overline{G_x'}$, then y is a type-1 node and the corresponding arc starting at $M(y)$ passes through x .

Proof. By definition, $M(a(y)) \subset x \subset M(y)$ and $M(b(y)) \subset M(y)$. Thus, y is a type-1 node. Obviously, the corresponding arc passes through x . \square

Claim 4. Let $y \in G_x$. If $a(y) \in G_x'$ and $b(y) \in G_x''$, then y is a type-2 node and the corresponding path starting at $M(y)$ passes through x .

Proof. Since $x \subseteq M(b(y))$ and $M(a(y)) \subset x$, then $y = a(y) \cup b(y) \subseteq x \cup M(b(y)) = M(b(y))$. By the definition, $M(y) = M(b(y))$. Thus, y is a type-2 node. Since $M(a(y)) \subset x \subseteq M(y)$, the corresponding path passes through x . \square

Let the numbers of type-1, type-2, type-3 nodes in G_x be n_1, n_2, n_3 respectively. By Claim 1 and Claim 2, the number of duplications in which x is mixed is $n_2 + n_3$ and equals the number of internal nodes in G_x . By Claim 3 and Claim 4, the number of paths passing x is $n_1 + n_2$ and equals the number of leaves of G_x . Since the subtree G_x is a binary, full, then $n_1 + n_2 = n_2 + n_3 + 1$. Thus, the fact is true. \square

Let duplication occur at p nodes of T ,

$$x_1, x_2, \dots, x_p,$$

and let there be t_i duplications at node x_i :

$$D_{i1}, D_{i2}, \dots, D_{it_i}.$$

Then, by Proposition 3.3,

$$\sum_{i=1}^p t_i = k_2 + k_3. \quad (2)$$

By Propositions 3.2, 3.4 and Formula (1), (2),

$$\begin{aligned}
c(T/G) &= \sum_{i=1}^p \sum_{j=1}^{t_i} c(D_{ij}) \\
&= \sum_{i=1}^p t_i + \sum_{1 \leq i \leq k_1} (|M(g_{1i})M(a(g_{1i}))| + |M(g_{1i})M(b(g_{1i}))|) + k_1 \\
&\quad + \sum_{1 \leq i \leq k_2} |M(g_{2i})M(a(g_{2i}))| + k_2 - (N - 1) + k_2 + k_3 \\
&= k_2 + k_3 + \sum_{1 \leq i \leq k_1} (|M(g_{1i})M(a(g_{1i}))| + |M(g_{1i})M(b(g_{1i}))|) \\
&\quad + \sum_{1 \leq i \leq k_2} |M(g_{2i})M(a(g_{2i}))| + k_2 \\
&= \sum_{1 \leq i \leq k_1} (|M(g_{1i})M(a(g_{1i}))| + |M(g_{1i})M(b(g_{1i}))|) \\
&\quad + \sum_{1 \leq i \leq k_2} (2 + |M(g_{2i})M(a(g_{2i}))|) + k_3 \\
&= c(G, T)
\end{aligned}$$

Hence, we have proved the following conjecture.

Theorem 3.1 (Mirkin-Muchnik-Smith Conjecture) *For any T and G with the same set of leaves, $c(T/G) = c(G, T)$.*

3.3 A linear-time algorithm

In this subsection, we shall present a linear-time algorithm for computing the cost function $c(T/G)$ for a gene tree G and a species tree T .

First, we compute the mapping M from G to T in linear time. Given a node $u \in G$, by the definition of mapping, its destination $M(u)$ is the lowest common ancestor of $M(a(u))$ and $M(b(u))$. This simple observation leads the following algorithm for computing the mapping from a gene tree to a species tree. In our algorithm, we define an auxiliary binary tree MT with information attached to various nodes. MT has the same structure as G . At each node $u \in MT$, we associated a pointer $m(u)$, which points to the destination $M(u)$ of u after it is computed. For $u, v \in G$, we use $lca(u, v)$ to denote the lowest common ancestor of u and v and $LCA(u, v)$ to denote the instruction for finding $lca(u, v)$. Finally, recall that in postorder, $i < j$ if and only if i is to the left of j or a descendant of j (see Aho *et al.*, 1974).

Algorithm 1

1. Generate a sequence, S , of LCA instructions by processing each node of G in postorder. For each internal node $u \in G$, generate an instruction: $LCA(M(a(u)), M(b(u)))$;
2. Execute the instruction sequence S . After finishing an instruction corresponding to $u \in G$, save the result to $MT(u)$.
3. Output $MT(u)$.

Theorem 3.2 *The mapping M can be computed in $O(n)$ time and $O(n)$ space on a RAM using Algorithm 1.*

Proof. Obviously, Step 1 takes $O(n)$ times and the instruction set S can be saved in $O(n)$ spaces.

In the instruction sequence S , the nodes involving in an instruction depend on the results of two previous instructions, and thus each instruction must be answered before processing the next. Therefore, we use the on-line algorithms for finding lowest common ancestors presented in Harel and Tarjan(1984) or Schieber and Vishkin (1988). Both algorithms takes $O(n)$ time and $O(n)$ space for executing the sequence S of $(n - 1)$ LCA instructions. But, the algorithm of Schieber and Vishkin is simple and its idea is as follows. The algorithm is based on the fact that on a line tree or a complete binary tree, each LCA query can be answered in $O(1)$ by direct calculation. In the preprocessing stage, the given tree T is divided into simple paths, and each path is mapped to a node in a complete binary tree B . Together with some addition information, a LCA query is answered by locating its path first and then the lca in the path. \square

Recall that for a node $u \in T$, the depth of u is the length of the simple path from the root to itself. The depth of the root is obviously 0. For computing the mapping cost function $c(G, T)$, we need to preprocess the tree T to get the depth $D(u)$ of each node $u \in T$ and then calculate the cost using information arrays D and MT .

Algorithm 2

1. Process each node of T in preorder. For each node u , calculate its depth $D(u)$.
2. Use Algorithm 1 to compute the mapping MT from G into T .
3. Compute the mapping cost c by process each internal node of MT .
Initially, $c = 0$. For each internal node $u \in MT(u)$, update the cost by

$$c = c + D(m(a(u))) + D(m(b(u))) - 2D(m(u)) - 2$$
 if $m(a(u)) \neq m(u)$ and $m(b(u)) \neq m(u)$,

$$c = c + D(m(a(u))) - D(m(u)) + 1$$
 if only $m(a(u)) \neq m(u)$, or

$$c = c + 1$$
 otherwise.

Theorem 3.3 *The mapping cost can be computed in $O(n)$ time and $O(n)$ space on a RAM using Algorithm 2.*

Proof. The correctness follows from the definition of the mapping cost. Observe that Both Step 1 and Step 3 take $O(n)$ time and $O(n)$ space. By Theorem 3.2, Step 2 takes $O(n)$ time and $O(n)$ space. Thus, Algorithm has the required time and space complexity. \square .

The previous algorithm for computing the mapping cost takes $O(n^3)$, which is due to Page(1994). By Theorem 3.1, the inconsistency cost is also computable in linear time and space on a RAM.

4 Conclusions

We have proved the Mirkin-Muchnik-Smith conjecture that the inconsistency function of a gene tree and a species tree equals the mutation cost function, and also provided a linear-time algorithm for computing the mutation cost function. These results reflect that the mutation cost function not only has biological meaningful, but also is very efficiently computable. Therefore, it is a good candidate in the future for measuring the similarity/dissimilarity between a gene tree and a species tree.

After this work was finished, the author found that Oliver Eulenstein and Martin Vingron also proved the conjecture independently (Eulenstein and Vingron, 1995).

Acknowledgement

The author would like to thank Ilya Muchnik for helpful discussions regarding to the topic in this paper through e-mail and Boris Mirkin for useful comments on the first draft of the paper. The author also thank the referee and Mike Waterman for suggestions and comments on revising the paper.

References

- [1] E.N. III Adams, N-trees as nestings: complexity, similarity and consensus, *Journal of Classification* 3(1986), 299-317.
- [2] A.V. Aho, J.E. Hopcroft and J.D. Ullman, *The Design and Analysis of Computer Algorithms*, Addison-Wesley, Reading, Mass., 1974.
- [3] J.P. Barthélemy, B. Leclerc, and B. Monjardet, On the use of ordered sets in problems of comparison and consensus of classifications, *Journal of Classification* 3(1986), 187-224.
- [4] B. DasGupta, X. He, T. Jiang, M. Li, J. Tromp and L. Zhang. On distance between phylogenetic trees. *Proc. of the 8th ACM-SIAM Symposium on Discrete Alg.*, 427-436, 1997.
- [5] O. Eulenstein and M. Vingron, On the Equivalence of Two Tree Mapping Measures, *Arbeitspapiere der GMD*, 936, Bonn, Germany.
- [6] W. M. Fitch, Distinguishing homologous and analogous proteins, *Syst. Zool.* 19(1970), 99-113.
- [7] M. Goodman, J. Czelusniak, G.W. Moore, A.E. Romero-Herrera and G. Matsuda, Fitting the gene lineage into its species lineage. A parsimony strategy illustrated by cladograms constructed from globin sequences, *Syst. Zool.* 28(1979), 132-163.
- [8] R. Guigó, I. Muchnik and T. Smith, Reconstruction of ancient molecular phylogeny, *Molecular Phylogenetics and Evolution* 6(1996), No. 2, 189-213.

- [9] D. Harel and R.E. Tarjan, Fast algorithms for finding nearest common ancestors, *SIAM J. on Comput.* 13(1984), No. 2, 338-355.
- [10] M.D. Hendy, C.H.C. Little and D. Penny, Comparing trees with pendant vertices labeled, *SIAM J. Appl. Math.* 44(1984), 1054-1067.
- [11] T. Margush and F. R. McMorris, Consensus n-Trees, *Bulletin of Mathematical Biology* 43(1981), 239-244.
- [12] B. Mirkin, I. Muchnik and T. Smith, A biologically meaningful model for comparing molecular phylogenies, *J. Comput. Biology* 2(1995), 493-507.
- [13] M. Nei, *Molecular Evolutionary Genetics*. Columbia University Press, New York.
- [14] J. E. Neigel and J. C. Avise, Phylogenetic relationship of mitochondrial DNA under various demographic models of speciation, *Evolutionary Processes and Theory*, 515-534, Academic Press, New York, 1986.
- [15] S. Ohno, *Evolution by gene duplication*. Springer-Verlag, Berlin, 1970.
- [16] R.D.M. Page, Maps between trees and cladistic analysis of historical associations among genes, organisms, and areas, *Syst. Biol.* 43(1994), 58-77.
- [17] P. Pamilo and M. Nei, Relationship between gene trees and species trees. *Mol. Bio. Evol.* 5 (1988), 568-583.
- [18] D. F. Robinson, Comparison of labeled trees with valency three, *Journal of Combinatorial Theory*, Series B, 11(1971), 105-119.
- [19] B. Schieber and U. Vishkin, On finding lowest common ancestors: simplification and parallelization, *SIAM J. on Comput.* 17(1988), No. 6, 1253-1262.
- [20] N. Takahata, Gene genealogy in three related population: Consistency probability between gene and population trees, *Genetics* 122(1989), 957-966.
- [21] M. S. Waterman and T. F. Smith, On the similarity of dendrograms, *Journal of Theoretical Biology* 73(1978), 789-800.
- [22] C.-I. Wu, Inference of species phylogeny in relation to segregation of ancient polymorphisms, *Genetics* 127(1991), 429-435.