



Multiple genome rearrangement: a general approach via the evolutionary genome graph

Dmitry Korkin* and Lev Goldfarb

Faculty of Computer Science, University of New Brunswick, Fredericton, NB, E3B 5A3, Canada

Received on January 24, 2002; revised and accepted on March 31, 2002

ABSTRACT

Motivation: In spite of a well-known fact that genome rearrangements are supposed to be viewed in the light of the evolutionary relationships within and between the species involved, no formal *underlying* framework based on the evolutionary considerations for treating the questions arising in the area has been proposed. If such an underlying framework is provided, all the basic questions in the area can be posed in a biologically more appropriate and useful form: e.g., the similarity between two genomes can then be computed via the nearest ancestor, rather than ‘directly’, ignoring the evolutionary connections.

Results: We outline an evolution-based general framework for answering questions related to the multiple genome rearrangement. In the proposed model, the evolutionary genome graph (EG-graph) encapsulates an evolutionary history of a genome family. For a set of all EG-graphs, we introduce a family of similarity measures, each defined via a fixed set of genome transformations. Given a set of genomes and restricting ourselves to the transpositions, an algorithm for constructing an EG-graph is presented. We also present the experimental results in the form of an EG-graph for a set of concrete genomes (for several species). This EG-graph turns out to be very close to the corresponding known phylogenetic tree.

Contact: dkorkin@unb.ca

Keywords: genome rearrangement; evolutionary model; gene transformations; longest common subsequence; ETS framework.

1 INTRODUCTION

The role of evolution in biology is so central that the two words are almost synonymous (Dobzhansky, 1973). As far as questions arising in connection with the multiple genome rearrangements are concerned, they have not yet been formulated and answered satisfactorily within an evolution-based formal framework. This is in spite of a well-known fact that genome rearrangements are

supposed to be treated mainly as an integral part of the evolutionary information, or relations, between the species.

In this paper we outline an evolution-based general framework for embedding some questions related to the multiple genome rearrangement. It should be quite clear that if the proposed model proves satisfactory, many additional kinds of useful questions could be addressed within such a model.

The main distinctive feature of the proposed model, as compared with the conventional genome rearrangement approaches, is that our model allows one to reconstruct evolutionary relationships between the species based on the evolutionary genome graph (EG-graph). The model also allows one to reconstruct other than binary relations between the species; for example, EG-graph can accommodate three genomes that have one closest common ancestor. Moreover, the model allows for a *natural* integration of point mutations into the overall scheme.

The proposed model was inspired by a much more general ‘evolution-based’ formalism for structural object representation within an inductive framework—evolving transformations system (ETS) model—developed over last 15 years (Goldfarb *et al.*, 2001; Goldfarb and Korkin, 2001). The ETS formalism was motivated by the integration of the concept of similarity into a generative inductive class representation. The resulting class representation can be constructed based on a small set of class examples and is conceptualized via a weighted set of transformations acting on the class progenitor (the common ancestor).

This paper is organized as follows. In Section 2, we present the basic definitions, including that of EG-graph, and formulate two versions of the basic problem of genome development (which are analogues of the corresponding genome rearrangement problem). In Section 3, we introduce an evolution-based model with reversals and insertions as non-local evolutionary mutations (transformations). We define the corresponding evolutionary genome graph and discuss its properties. Next, we discuss the relationship between evolutionary transformations and the traditional genome rearrangement transformations.

*To whom the correspondence should be addressed.

We also introduce an evolutionary similarity measure for this model. At the end of this section, we introduce context-sensitive genome transformations and discuss how they could make the above model more accurate. In Section 4, we outline a method of constructing the EG-graph for a given set of genomes. We present the algorithm and analyse its complexity. In Section 5, we present the results of an experiment, i.e., an EG-graph for genomes of several species, and compare them to the corresponding known phylogenetic trees. In the last section, we conclude with a brief summary of the paper and propose some future directions.

2 BASIC DEFINITIONS AND PROBLEM FORMULATION

Some of the following concepts are standard, while the problems are formulated for the first time and rely on the evolutionary genome graph.

2.1 Basic definitions

DEFINITION 1. Let $\Sigma_0 = \{A, C, G, T\}$ be an alphabet. A **gene** is defined to be an element of the set of strings $S = \Sigma_0^* \setminus \{\Lambda\}$, and a **genome G** is defined to be a sequence of genes: $G = (s_1, s_2, \dots, s_n)$, $s_k \in S$, $1 \leq k \leq n$. The corresponding string $G = s_1 \circ s_2 \circ \dots \circ s_n$ is called the **genome sequence**. The **length**, $|G|$, of a genome is defined to be the length of its genome sequence. The set of all genomes is denoted as Γ .

Note that the circular genomes can and will also be represented as sequences (see Section 4.1). The next definition formally introduces the concept of gene order transformation, which leaves the length of the genome unchanged.

DEFINITION 2. Let $A = (a_1, a_2, \dots, a_n)$ and $B = (b_1, b_2, \dots, b_n)$ be any genomes of the same length, such that there exists a permutation $\rho : \{1, 2, \dots, n\} \rightarrow \{1, 2, \dots, n\}$ satisfying $b_i = a_{\rho(i)}$ or $b_i = a'_{\rho(i)}$, where a' is a reversed a . A mapping $f : \Gamma \rightarrow \Gamma$, $f(A) = B$, is called a **gene order transformation** for genome A , or simply **gor-transformation**. The set of all gor-transformations is denoted as Φ_{gor} .

The following definition formally introduces the concept of genome rearrangement.

DEFINITION 3. A **rearrangement of genome G** terminating in genome G_m is a sequence (f_1, f_2, \dots, f_m) , $m \geq 1$, $f_i \in \Phi_{gor}$, defined inductively as follows:

1. f_1 is a gor-transformation for genome G and $G_1 = f_1(G)$
2. $\forall i, 1 < i \leq m$, f_i is a gor-transformation for genome G_{i-1} and $G_i = f_i(G_{i-1})$.

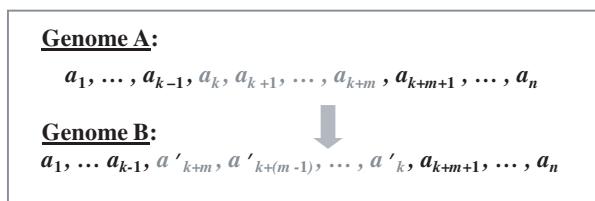


Fig. 1. A reversal of subsequence $A_1 = a_k, a_{k+1}, \dots, a_{k+m}$.

The next two definitions review the concepts of subsequence and the multiple longest common subsequence.

DEFINITION 4. Let a be a sequence of length n over another alphabet $\Sigma = \{a, b, \dots, z\}$. If sequence $a = s_1 s_2 \dots s_n$, $s_i \in \Sigma$, then sequence $b = s_{i_1} s_{i_2} \dots s_{i_k}$ is called a **subsequence** of a , if $\forall j, 1 \leq j \leq k$:

$$1 \leq i_j \leq n,$$

and for all s and t , $1 \leq s < t \leq k : i_s < i_t$, where k is the length of b .

DEFINITION 5. Let $S = \{a_1, a_2, \dots, a_d\}$ be a set of sequences over alphabet $\Sigma = \{a, b, \dots, z\}$, of lengths n_1, n_2, \dots, n_d , correspondingly. The **multiple longest common subsequence (MLCS)** for a set S is a sequence b such that:

1. b is a subsequence of a_i , $\forall i$;
2. b is the longest sequence satisfying 1.

In case $d = 2$, MLCS is called simply the *longest common subsequence (LCS)*.

There are two particular kinds of gene order transformations, i.e., reversals and transpositions, that are of interest.

DEFINITION 6. Let $A = (a_1, a_2, \dots, a_n)$ and $B = (b_1, b_2, \dots, b_n)$ be genomes and let f be a gor-transformation for genome A such that $B = f(A)$. The gor-transformation f is called a **reversal** of a subsequence $A_1 = (a_k, a_{k+1}, \dots, a_{k+m})$ of genes, if the corresponding permutation ρ (see Definition 2) is such that there exist integers $k, m \geq 1$ satisfying (see Figure 1):

$$b_i = a_i, 1 \leq i \leq k-1 \quad \text{or} \quad k+m+1 \leq i \leq n \quad (1)$$

$$b_{k+i} = a'_{k+(m-i)}, 0 \leq i \leq m. \quad (2)$$

DEFINITION 7. Let $A = (a_1, a_2, \dots, a_n)$ and $B = (b_1, b_2, \dots, b_n)$ be two genomes and let f be a gor-transformation for genome A such that $B = f(A)$. The gor-transformation f is called a **transposition** of a subsequence $A_1 = (a_k, a_{k+1}, \dots, a_{k+m})$ of genes, if the corresponding permutation ρ is such that there exist $j, k, m \geq 1$ satisfying (see Figure 2):

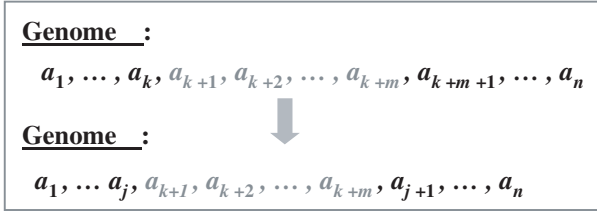


Fig. 2. A transposition of $A_1 = a_k, a_{k+1}, \dots, a_{k+m}$.

- (1) $b_i = a_i, 1 \leq i \leq \min(k, j) \text{ or } \max(k, j) + 1 + m \leq i \leq n$
- (2) $b_{j+i} = a_{k+i}, 1 \leq i \leq m$
- (3) $b_{i+m} = a_i, j + 1 \leq i \leq k$

(note that we assume $j \leq k$; the case $j > k$ can be defined similarly).

The next definition formally introduces the concept of insertion transformation, which changes the length of the genome.

DEFINITION 8. Let $A = (a_1, a_2, \dots, a_n)$ and $B = (b_1, b_2, \dots, b_{n+m}), m > 0$, be two genomes. A mapping $f : \Gamma \rightarrow \Gamma$, where $f(A) = B$, is called **insertion transformation** (for a genome A) of sequence $A_1 = (c_1, c_2, \dots, c_m)$ of genes, or simply **ins-transformation**, if there exists $k, 1 \leq k \leq n$, such that:

- (1) $b_i = a_i, 1 \leq i \leq k$
- (2) $b_{k+i} = c_i, 1 \leq i \leq m$
- (3) $b_{i+k+m} = a_{i+k}, 1 \leq i \leq n - k$.

The set of all ins-transformations will be denoted as Φ_{ins} . The set of all gor- and ins-transformations will be denoted as Φ .

One should note the use of pronouns ‘for’ and ‘of’ in the above definitions (and below).

In the next definition, we introduce the evolutionary version of the concept of genome rearrangement by adding to the set of gor-transformations the set of ins-transformations.

DEFINITION 9. A **development of genome G** terminating in genome G_m is a sequence $(f_1, f_2, \dots, f_m), m \geq 1, f_i \in \Phi$, defined inductively as follows:

- 1. f_1 is a gor- or ins-transformation for genome G and $G_1 = f_1(G)$;
- 2. $\forall i, 1 < i \leq m, f_i$ is a gor- or ins-transformation for genome G_{i-1} and $G_i = f_i(G_{i-1})$.

We will also need the concept of closest common ancestor useful from the evolutionary point of view.

DEFINITION 10. A genome C is called a **common ancestor** for a set of genomes Γ_1 , if $\forall G \in \Gamma_1$ either $G = C$ or there exists a development of genome C terminating in G .

Let Γ_{1A} be the set of all common ancestors of Γ_1 . A genome C is called a **closest common ancestor** for Γ_1 , if it belongs to set $\Gamma_{1C}, \Gamma_{1C} \subseteq \Gamma_{1A}$, such that for any common ancestor G (of Γ_1) not in Γ_{1C} there exists $C' \in \Gamma_{1C}$ and a development of G terminating in C .

2.2 Evolutionary genome graph and problem formulation

The current evolutionary approaches to genome rearrangement have roots in phylogenetic approaches to the study of protein sequences and have appeared during the last ten years (Hannenhalli *et al.*, 1995; Sankoff *et al.*, 1992, 1996; Sankoff and Blanchette, 1998, 1999; El-Mabrouk, 2001).

In our approach, the evolutionary genome graph is supposed to represent the evolutionary dependencies for the given set of genomes (see Figure 3).

DEFINITION 11. In a directed graph $DG, DG = (V, E)$, a vertex from V all edges of which are outgoing will be called a **source node**. A vertex from V all edges of which are incoming will be called a **sink node**.

DEFINITION 12. Given a finite set of genomes $\Gamma_1 \subseteq \Gamma$, the **evolutionary genome graph**, or simply EG-graph, for Γ_1 is defined as connected directed labeled graph $DG, DG = (V, E), l_V : V \rightarrow \Gamma, l_E : E \rightarrow \Phi$ (See Definition 8), such that:

- (1) there exists only one source node, $V_0, V_0 \in V$
- (2) l_V is an injective mapping and $\Gamma_1 \subseteq l_V(V)$
- (3) if $G_0 = l_V(V_0)$, then G_0 is a closest common ancestor for Γ_1 .

It is easy to see that, in general, EG-graph may not be a tree: there could be two or more different paths between two vertices, e.g., one path obtained by performing the insertion of a gene g_1 for a genome G followed by a reversal of another gene g_2 for G and another path obtained by performing the reversal of g_2 first followed by the insertion of g_1 for the genome G .

Having introduced the concept of EG-graph, we can formulate the basic problem of genome development (which is an analogue of the corresponding genome rearrangement problem).

PROBLEM 1. Given a finite set of genomes $\Gamma_1 = \{G_1, G_2, \dots, G_n\}$, construct an EG-graph, $DG, DG = (V, E)$, such that $\Gamma_1 \subseteq l_V(V)$.

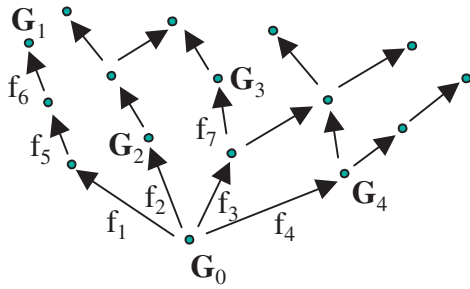


Fig. 3. An evolutionary genome graph.

Next, given an EG-graph, $DG, DG = (V, E)$, in order to be able to compare a pair of genomes with respect to the set of transformations $l_E(E)$, one needs to specify a similarity measure

$$\mu_{DG} : l_V(V) \times l_V(V) \rightarrow \mathbb{R}^+.$$

In *evolutionary* approaches, one of the basic and most natural ways to introduce such a measure is by means of the evolutionary paths, i.e., by means of the sequence of transformations to the two genomes from their closest common ancestor (Figure 4). Since the set $l_E(E)$ of transformations may vary, depending on the EG-graph, we obtain the family of similarity measures:

$$M = \{\mu_{DG} | DG = (V, E) \text{ is an EG-graph}\}.$$

Having specified a similarity measure on a set of genomes, we can formulate an important, although computationally more complex, problem (which is an analogue of the corresponding optimal genome rearrangement problem).

PROBLEM 2. *Given a finite set of genomes, $\Gamma_1 = \{G_1, G_2, \dots, G_n\}$, and a family of similarity measures, $M = \{\mu_{DG}\}$, construct an EG-graph, $DG, DG = (V, E)$, such that $\Gamma_1 \subseteq l_V(V)$ and at which the minimum of the following function (the mean value) is achieved over all EG-graphs with $\Gamma_1 \subseteq l_V(V)$:*

$$\rho_{DG}(\Gamma_1) = \sum_{G_i, G_j \in \Gamma_1} \mu_{DG}(G_i, G_j) / n.$$

3 A NOVEL GENOME DEVELOPMENT MODEL

In this section, we introduce a novel evolution-based framework for genome development. First, we consider the case of two basic types of transformations, reversals and ins-transformations, and discuss their relationships with gen-transformations. Then, we introduce a similarity

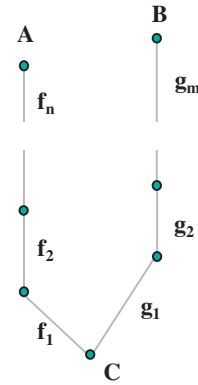


Fig. 4. The closest common ancestor of two genomes, A and B

measure based on these basic types of transformations. Finally, we briefly discuss a new concept—the context of transformation—and show how it can be encapsulated in the model.

3.1 An evolution-based genome development model

Consider a set of genes $S = \{g_1, g_2, \dots, g_n\}$, i.e., a set of strings over the alphabet $\Sigma_0 = \{A, C, G, T\}$. We now specify some biologically reasonable restrictions we impose on the above set Γ of genomes (Definition 1) and EG-graphs (Definition 12). Set Γ of genomes is now restricted to those genomes that do not have repeated genes:

$$\Gamma = \{G | G = (g_{i_1}, g_{i_2}, \dots, g_{i_k}) \text{ and } i_1, i_2, \dots, i_k \in \{1, 2, \dots, n\}, i_1 \neq i_2 \neq \dots \neq i_k\}.$$

We also specify a particular class of the corresponding EG-graphs (Definition 12). Given a finite set of genomes $\Gamma_1 \subseteq \Gamma$, the restrictions on an EG-graph $DG, DG = (V, E)$, for Γ_1 are specified as follows:

- (1) $\forall f \in l_E(E)$, f is either a reversal or an insertion
- (2) $\forall G \in \Gamma_1, \forall$ genome development of G_0 terminating in $G, \forall g \in S$, such that g is present in genome G , g cannot be reversed twice (back to original form) by any of the reversals from the genome development
- (3) if V_S is the set of sink nodes (Def. 11) of DG , then $l_V(V_S) = \Gamma_1$.

Note that the additional restrictions on the EG-graph, in particular, allow us to avoid the case where a reversal transformation can be applied infinitely many times for a genome. Moreover, the following lemma is valid:

LEMMA 1. *With the above restrictions, the connected directed labeled graph $DG = (V, E)$ is also an acyclic graph.*

The above Lemma implies two related and very important properties: first, the genome can never return to one of its previous states, and, as a consequence of this, second, if one allows copying (that is, insertion) of each of the genes only once (in a more general model, several times), then the process of construction of all possible genomes is finite.

The model specified above has strong connections with the traditional genome transformations. Namely, it is not difficult to obtain the following results. (Henceforth, when necessary, we will denote a transformation f of a sequence of genes X as f_X .)

LEMMA 2. Let A and B be two genomes, $A, B \in \Gamma$, and let genome C be their closest common ancestor. Then:

- (1) if $B = f_X(A)$, where f_X is a transposition of a sequence X of genes, then there exist two insertions of this sequence of genes f'_X and f''_X , such that $A = f'_X(C)$, $B = f''_X(C)$;
- (2) if $B = f(A)$, where f is a reversal, then $A = C$.

In other words, those genomes that are related to each other in a traditional model by either a reversal or a transposition, are also related in our model, but now via the closest ancestor and the corresponding gor- or ins-transformations. This partly explains the need for an evolutionary similarity measure for genomes to be based on the transformations leading to these genomes.

Next, we introduce a similarity measure for the above model. Before specifying the similarity measure formally, we want to assign a weight to each of the reversal and insertion transformations used in the corresponding EG-graph. There are many ways to assign a weight to a genome transformation. One of the standard ways is considering it with the size of the inserted, reversed, or transposed substring. In our model, we will use the following weighting scheme.

DEFINITION 13. For any reversal or insertion f_A of a sequence of genes A , the **weight** of this transformation is defined as $w(f_A) = |A|$, where $|A|$ is the length of A (see Definition 1).

Having assigned the weights to transformations, we can choose the following similarity measure.

DEFINITION 14. Given an EG-graph DG , let A and B be two genomes, $A, B \in \Gamma$, and let genome C be their closest common ancestor. Moreover, let $A = f_n \circ f_{n-1} \circ \dots \circ f_1(C)$, $B = g_m \circ g_{m-1} \circ \dots \circ g_1(C)$ (see Figure 4). Then, the **similarity measure** is defined as

$$\mu_{DG}(A, B) = \sum_{i=1}^n w(f_i) + \sum_{i=1}^m w(g_i).$$

Thus, having specified the class of EG-graphs, weighting scheme, and the similarity measure, it becomes possible to introduce the relationships among the genomes. Namely, given a set of genomes Γ_1 , we can, first, construct a particular EG-graph, thus specifying the particular mutation pathways that lead to each of the genomes in Γ_1 . Then, any genome in Γ_1 can be compared with any other genome in Γ_1 and their similarity measure, in terms of similar transformations occurring in the closest common ancestor of the above genomes, can be calculated.

3.2 The case of context-sensitive gene transformations

In this section, we discuss the concept of context-sensitive genome transformations. Why is such a concept useful? It is natural to assume that the process of a genome rearrangement (and development) should depend on the structure of genome being rearranged. In other words, the transformation for a genome may depend on a particular region of this genome (the case of a *local context*), or it may depend on the regions that are not close to each other and even on the entire genome (the case of a *global context*). Since the concept of a global context-sensitive transformation is very complex, in this introductory paper, we will discuss the local context-sensitive transformations only. We next consider the concept of the local context for both types of transformations, reversals and insertions.

DEFINITION 15. Given a set of genomes, Γ_1 , and a transformation f_X , where f_X is either reversal or ins-transformation, the **local context** of f_X is a pair of sequences (C_1, C_2) , $C_1, C_2 \in \Sigma^*$, such that

- (1) if f_X is a reversal and X' is the reversed X , then $\forall A, B \in \Gamma_1$, where $B = f_X(A)$,

$$A = A^1 \circ C_1 \circ X \circ C_2 \circ A^2$$

and

$$B = B^1 \circ C_1 \circ X' \circ C_2 \circ B^2,$$

where $A^1, A^2, B^1, B^2 \in \Sigma^*$;

- (2) if f_X is an ins-transformation, then $\forall A, B \in \Gamma_1$, where $B = f_X(A)$,

$$A = A^1 \circ C_1 \circ C_2 \circ A^2$$

and

$$B = B^1 \circ C_1 \circ X \circ C_2 \circ B^2,$$

where $A^1, A^2, B^1, B^2 \in \Sigma^*$.

Note that C_1 and/or C_2 can be the null string. In the case when both C_1 and C_2 are null strings, the corresponding transformation f_X is said to be **context-free**.

How does the introduction of context affect the similarity measure? There are different ways to redefine the similarity measure to reflect the presence of a context. It goes without saying that the context-sensitive transformation, when applied to some genome, acts more discriminatively than the ‘same’ transformation but without any context. Therefore, after the application of a context-free operation f_X , the resulting genome $B = f_X(A)$ might be considered *farther* away from the original genome A than the one obtained by application of the ‘same’ transformation f_X but with non-empty context. In other words, in view of the additive nature of the similarity measure, the weight of a context-free transformation should be larger than that of the same transformation but with a non-empty context. Below, we give one of the possible ways to define such weighting scheme for context-sensitive transformations.

DEFINITION 16. Let f_X be a transformation and (A, B) be its context. Then, the **context-sensitive** weight of transformation f_X is defined as

$$w_{CS}(f_X) = w_{CF} \frac{|X|}{|X| + |A| + |B|},$$

where w_{CF} is a ‘standard’ weighting scheme for a (context-free) transformation introduced in Definition 13.

The weighting scheme in Def. 16 has two important features. First of all, for any context-free transformation f_X , its context-sensitive weight, $w_{CS}(f_X)$, is equal to the ‘standard’ weight of f_X , and the bigger the context of f_X , the smaller its context-sensitive weight. The latter, as was already mentioned above, can be explained by the fact that the bigger the context of a transformation, the more specific this transformation (when applied to a genome), and thus, the resulting new genome should be closer to its ancestor than the one obtained by applying the context-free version of the same transformation f_X . Second, when solving Problem 2, defined above, a smaller weight of a context-sensitive transformation makes it preferable, since the similarity measure between two genomes based on context-sensitive transformations will also be smaller in comparison with the same measure based on the same transformations but without the contexts.

4 IMPLEMENTATION: GENOME REARRANGEMENT USING TRANSPOSITIONS ONLY

In this section, we discuss how to reconstruct an EG-graph, given a set of genomes. We consider a basic type of rearrangement that uses only transpositions. The EG-graph reconstruction algorithm, described in this section, uses the idea of multiple longest common subsequence (see Definition 5).

4.1 Some basic ideas and assumptions

The basic, and computationally the simplest, is the case when we restrict ourselves to the set of genomes Δ , $\Delta \subseteq \Gamma$, possessing the property that any of its genomes can be obtained from some other of its genomes by applying a finite set of transpositions of genes from S . In an evolution-based model, this means that the same set of genes is to be inserted (possibly in a different order) by the corresponding ins-transformations into the protogenome. This assumption results in the following lemmas:

LEMMA 3. *Each of genomes from Δ consists of the same number of genes from S .*

LEMMA 4. *Given a finite set of genomes $\Gamma_1 \subseteq \Delta$, let DG , $DG = (V, E)$, be any EG-graph for Γ_1 (see Section 3.1). Then, for any set $V_0 \subseteq V$, the MLCS of V_0 will contain at least one gene from S .*

The last lemma follows from the previous one and allows us to rely on the concept of MLCS in the algorithm discussed next.

4.2 Algorithm

The algorithm outputs one of the possible EG-graphs, $DG = (E, V)$, given a set of genomes Γ_1 composed of genes from a set S . As a preprocessing stage, to obtain a string representation for each of the circular genomes, i.e., to find a common gene along which each of the given circular genome will be cut, the following simple algorithm is applied. For each of the genes, all of the (circular) genomes are cut along this gene and for the resulting strings the MLCS is constructed. The final cut is determined by choosing genome cuts associated with a gene for which the MLCS found is the longest one (ties are broken arbitrarily). The algorithm consists of two parts. In the first part, the common ancestor of all genomes in Γ_1 is constructed. In the second (main) part, the transformations leading to each of the genomes are consecutively extracted. The pseudocode of the algorithm is presented below (`// ... //` marks the comments).

Algorithm EG-graph

Input: Γ_1 , $|\Gamma_1| = N$

Output: $DG = (V, E)$, where $V = \{v_1, v_2, \dots, v_K\}$ is the set of labeled vertices and $E = \{e_1, e_2, \dots, e_M\}$ in the set of labeled edges.

```
// 1. Construct ancestor //
Ancestor =  $v_0$  = MLCS( $\Gamma_1$ )
// 2. Construct transformations //
cur_level =  $V = \{v_0\}$ ;
level = 0; cur_G $_v$ ( $v_0$ ) =  $\{1, 2, \dots, N\}$ ;
// contains indices of all the genomes //
While cur_level is not empty do
```

```

{ For all vertices  $v$  in cur_level
{New_level =  $\emptyset$ ;
// cur_Gv: indices of all current genomes correspond-
ing to a current vertex from Cur_level //
cur_S = MLCS(cur_Gv);
cur_Tr = Transform(S);
// Transforms the MLCS cur_S to the sequence of
ins-transformations, cur_Tr //
For each transformation  $f$  in cur_Tr
{(V1, E1, cur_Gv) = Children_v(v, f, cur_Gv, V, E);
// Returns new vertices (V1) together with edges
(E1), connecting vertex  $v$  and new vertices, based
on the transformation  $f$  and the set of genomes
corresponding to vertex  $v$ , cur_Gv, plus the set of in-
dices Cur_Gv( $w$ ) corresponding to genomes that are
children of the genome corresponding to a vertex
 $w$  in V1//
V = V  $\cup$  V1; E = E  $\cup$  E1;
} // End For //
New_level = New_level  $\cup$  V1;
// Adds to a new level the vertices, corresponding to the
last transformation in the sequence cur_Tr//
} //End For all vertices  $v$  //
cur_level = new_level;
level = level + 1;
} //End While//
END //Algorithm//

```

The algorithm's pseudocode presented above depends on two basic subroutines. While the first subroutine, MLCS(), is discussed in Hakata and Imai (1998), the second one, Children_v (), can be described as follows. Given: a vertex v_A , corresponding to a genome A , the set of indices $cur_G_v(v)$ of those genomes in Γ_1 which have A as a common ancestor, and an ins-transformation f ; the subroutine Children_v () constructs

- (1) the set of vertices $V_1 = \{v_{A1}, v_{A2}, \dots, v_{Ad}\}$ corresponding to the set $\{A_1, A_2, \dots, A_d\}$ of immediate descendants of A obtained by applying f to A ;
- (2) the set of edges E_1 , each of which is labeled by f , connecting each vertex in V_1 with v_A ;
- (3) for each of A_i —the set of indices $cur_G_v(A_i)$ corresponding to those genomes in Γ_1 which have A_i as a common ancestor.

4.3 Example

This example is presented for illustrative purposes only. However, it captures all the necessary aspects of the algorithm as well as allowing to comparison of the results of 'traditional' genome rearrangement with those for the evolutionary genome rearrangement represented by the corresponding EG-graph. Figure 5 and Figure 6 present the initial data and the resulting EG-graph, respectively.

$G_1 = D A B E C F G H$;
 $G_2 = A D G B C F H E$;
 $G_3 = A G B D C F E H$;
 $G_4 = G A B D C F E H$;
 $G_5 = E G C A F D B H$;
 $G_6 = C E A F B G H D$.

Fig. 5. The input data, Γ_1 , consisting of 6 genomes composed of the following genes $\{A, B, \dots, H\}$.

4.4 Complexity

To estimate the computational time complexity of the algorithm, we use the following result.

LEMMA 5. Suppose $|\Gamma_1| = d$, $|S| = s$. Then:

1. The algorithm traverses the EG-graph only once;
2. For each level, the MLCS algorithm is performed no more than $d/2$ times;
3. There are no more than s levels in an EG-graph;
4. Transform () is $O(dL)$, where

$$L = \max\{|\mathbf{G}_1|, |\mathbf{G}_2|, \dots, |\mathbf{G}_d|\};$$

5. In the main algorithm, Children () is $O(ds)$.

Based on the above Lemma, it is not difficult to estimate the time complexity of the EG-graph algorithm. Namely, let $O(T)$ be the time complexity for MLCS subroutine. Then the following result is true.

THEOREM 1. The time complexity of the algorithm EG-graph is $O(sd(T + L))$.

Finally, we note that the following estimation of the time complexity of the MLCS algorithm, based on the dominant point approach (Hakata and Imai, 1998), points to a substantial advantage of the latter approach as compared to the known dynamic programming approaches for solving this problem.

THEOREM 2 (HAKATA AND IMAI, 1998). The MLCS problem for $d(d \geq 3)$ strings of length n can be solved in time

$$O(nsd + |D|sd(\log^{d-3} n + \log^{d-2} s)),$$

where $|D|$ is the size of the set of all dominant positions.

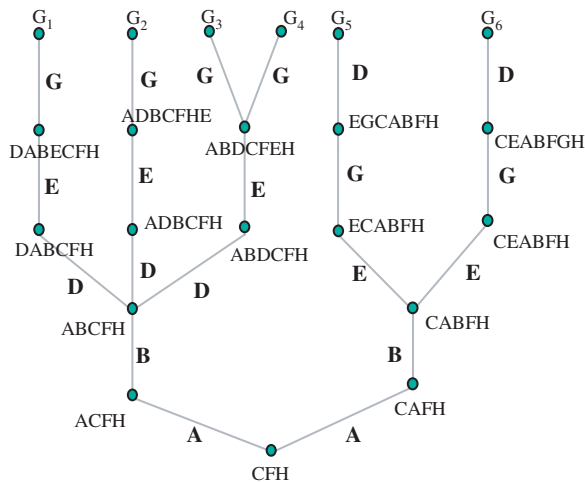


Fig. 6. An EG-graph for the set of genomes Γ_1 in Figure 5.

N	Organism	Group	Our genome cod
1	Albinaria coerulea	MOLuscs	(FAEOILKCDJBHG)
2	Asterina pectinifera	ECHinoderms	(JCDFGOBAHEILK)
3	Drosophila yakuba	ARTthropod	(BHILKJCFDEGOA)
4	Homo sapiens	CHOrdat	(ABHILKJCEDFGO)
5	Katharina tunicata	MOLuscs	(HIKLFDEOGAJCB)
6	Lumbricus terrestris	ANNeli	(HILJGOKFEDACB)

Our abbreviations for the genes

ND1	A	ND4	D	ND6	G	COX3	J	CYT B	O
ND2	B	ND4L	E	COX1	H	ATP6	K		
ND3	C	ND5	F	COX2	I	ATP8	L		

Fig. 7. Mitochondrial genomes and their assumed monophyletic groupings used in our experiments (adapted from (Blanchette *et al.*, 1999), Table 1).

5 EXPERIMENTAL RESULTS

Selected results of our experiments are shown in Figures 7–9. Note that under the assumed constraints the corresponding EG-graph must be a tree. The resemblance of our tree with those shown in Figure 8 is quite apparent. However, instead of being an unrooted tree, the resulting EG-graph is always a rooted tree whose root is a common ancestor of all six genomes and encapsulates all their common conserved parts.

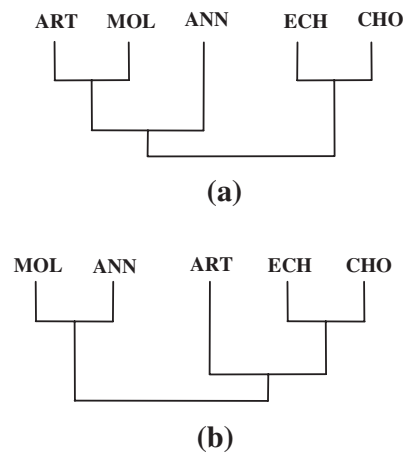


Fig. 8. Two alternative evolutionary trees for the species in Figure 7(a) currently most widely accepted view (b) a minimal breakpoint tree (adapted from Blanchette *et al.* (1999, Figures 1, 4)).

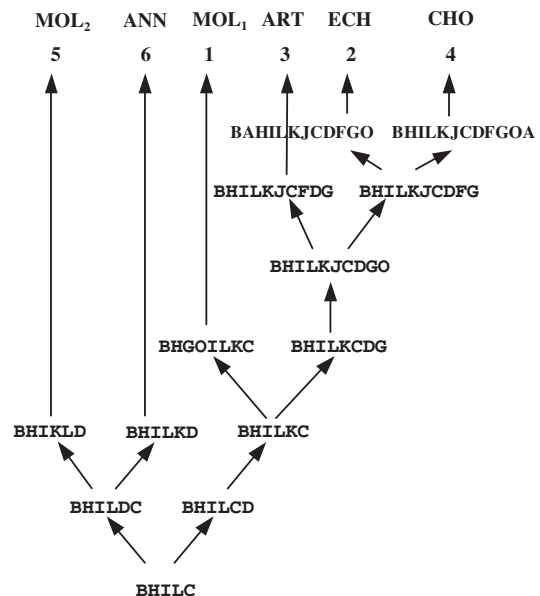


Fig. 9. The constructed evolutionary genome graph for the species in Figure 7.

6 DISCUSSIONS AND FUTURE RESEARCH

In this paper, we introduced an alternative, evolution-based, approach to the study of genome rearrangements. Within the approach, the development of genomes can be represented by a special directed graph, called an EG-graph, with labeled edges and vertices. One of the main advantages of the evolution-based approach is that it allows one to reconstruct other than just binary relations among the genomes: e.g., an EG-graph can represent three genomes that have one closest common ancestor. In order

to represent such relationships, we relied on the idea of a multiple longest common subsequence. Although the classical dynamic programming methods for computing MLCS can be used, practically, for the case of only two or three strings, the dominant-points based method allows one to obtain a MLCS for a much larger set of genome sequences.

As far as the future research directions are concerned, there are many that one can choose to follow. First of all, as the next step, one can consider the reconstruction of an EG-graph for more complex models of genome rearrangement, e.g., models based on the following genome transformations:

- (1) *transpositions and reversals of genes*
- (2) *transpositions and insertions of genes*

These models are not obvious and need careful study.

Next, one can consider the computationally more complex problem of reconstructing an optimal EG-graph (see Problem 2 in Section 2.B).

Another direction is the reconstruction of the context-sensitive transformations. One of the possible approaches, when the optimal EG-graph is not necessary, is to search for the context of transformations on the basis of the context-free EG-graph. Finally, one can consider a model of genome rearrangement with the presence of noise, i.e., point mutations. The weighting scheme should take this fact into consideration in such a way that the presence of some point mutations affects the similarity measure between the two genome sequences.

ACKNOWLEDGEMENTS

We would like to thank Patricia Evans for suggesting to the first author the genome rearrangement problem as an

interesting topic (for the application of the new model) as well as the source for the data.

REFERENCES

- Blanchette,M., Kunisawa,T. and Sankoff,D. (1999) Gene order breakpoint evidence in animal mitochondrial phylogeny. *J. Mol. Evol.*, **49**, 193–203.
- Dobzhansky,T. (1973) Nothing in biology makes sense except in the light of evolution. *The American Biology Teacher*, **35**, 125–129.
- El-Mabrouk,N. (2001) Reconstructing an ancestral genome using minimum segments duplications and reversals. *J. Comput. Syst. Sci.*, special issue on computational molecular biology, to appear.
- Goldfarb,L., Golubitsky,O. and Korkin,D. (2001) What is a structural representation? Submitted to the special issue of *Pattern Recognition* on the ETS model. In Goldfarb,L. (ed.),
- Goldfarb,L. and Korkin,D. (2001) Molecular classes and the structural representation: Towards theoretical foundations and unification of cheminformatics and chemistry. Submitted to the special issue of *Pattern Recognition* on the ETS model.
- Hakata,K. and Imai,H. (1998) Algorithms for the longest common subsequence problem for multiple strings based on geometric maxima. *Optimiz. Meth. Software*, **10**, 233–260.
- Hannenhalli,S, Chappey,C., Koonin,E. and Pevzner,P. (1995) Genome sequence comparison and scenarios for gene rearrangements: a test case. *Genomics*, **30**, 299–311.
- Sankoff,D. and Blanchette,M. (1998) Multiple genome rearrangement and breakpoint phylogeny. *J. Comput. Biol.*, **5**, 555–570.
- Sankoff,D. and Blanchette,M. (1999) Probability models for genome rearrangement and linear invariants for phylogenetic inference. *RECOMB*, 302–309.
- Sankoff,D., Leduc,G., Antoine,N., Paquin,B., Land,B.F. and Cedergren,R. (1992) Gene order comparisons for phylogenetic inference: Evolution of the mitochondrial genome. *Proc. Natl Acad. Sci. USA*, **89**, 6575–6579.
- Sankoff,D., Sundaram,G. and J.Kececioglu (1996) Steiner points in the space of genome rearrangements. *Int. J. Foundations Comput. Sci.*, **7**, 1–9.