

Biology Computing

Bhaskar DasGupta
Department of Computer Science
Rutgers University
Camden, NJ 08102, USA
E-mail: bhaskar@crab.rutgers.edu

Lusheng Wang *
Dept. of Computer Science
City University of Hong Kong
Tat Chee Avenue
Kowloon, Hong Kong
E-mail: lwang@cs.cityu.edu.hk

July 23, 1999

Keywords: Special purpose computing related to biology and biotechnology; DNA sequencing; evolutionary trees: construction and comparison; multiple sequence alignment problems.

1 Introduction

The modern era of molecular biology began with the discovery of the double helical structure of DNA. Today, sequencing nucleic acids, the determination of genetic information at the most fundamental level, is a major tool of biological research [79]. This revolution in biology has created a huge amount of data at great speed by directly reading DNA sequences. The growth rate of data volume is exponential. For instance, the volume of DNA and protein sequence data is currently doubling every 22 months [55]. One important reason for this exceptional growth rate of biological data is the medical use of such information in the design of diagnostics and therapeutics [24, 50]. For example, identification of genetic markers in DNA sequences would provide

*Supported in part by Hong Kong Research Council.

important informations regarding which portions of the DNA are significant, and would allow the researchers to find many disease genes of interest (by recognizing them from the pattern of inheritance). Naturally, the large amount of available data poses a serious challenge in storing, retrieving and analyzing biological information.

A rapidly developing area, *computational biology*, is emerging to meet the rapidly increasing computational need. It consists of many important areas such as information storage, sequence analysis, evolutionary tree construction, protein structure prediction, and so on [24, 50]. It is playing an important role in some biological research. For example, sequence comparison is one of the most important methodological issues and most active research areas in current *biological sequence analysis*. Without the help of computers, it is almost impossible to compare two or more biological sequences (typically, at least a few hundred character long).

In this chapter, we survey recent results on evolutionary tree construction and comparison, computing syntenic distances between multi-chromosome genomes, and multiple sequence alignment problems.

Evolutionary trees model the evolutionary histories of input data such as a set of species or molecular sequences. Evolutionary trees are useful for a variety of reasons, for example, in homology modeling of (DNA and protein) sequences for diagnostic or therapeutic design, as an aid for devising classifications of organisms, in evaluating alternative hypotheses of adaption and ancient geographical relationships (for example, see [25, 37] for discussions on the last two applications). Quite a few methods are known to construct evolutionary trees from the large volume of input data. We will discuss some of these methods in this chapter. We will also discuss methods for comparing and contrasting evolutionary trees constructed by various methods to find their similarities or dissimilarities, which is of vital importance in computational biology.

Syntenic distance are a measure of distance between multi-chromosome genomes (where each chromosome is viewed as a set of genes). Applications of computing distances between genomes can be traced back to the well-known **Human Genome Project**, whose objective is to decode this entire DNA sequence and to find the location and ordering of genetic markers along the length of the chromosome. These genetic markers can be used, for example, to trace the inheritance of chromosomes in families and thereby to find the location of disease genes. Genetic markers can be found by finding DNA polymorphisms, i.e., locations where two DNA sequences “spell” differently. A key step in finding DNA polymorphisms is the calculation of

the *genetic distance*, which is a measure of the correlation (or similarity) between two genomes.

Multiple sequence alignment is an important tool for sequence analysis. It can help extracting and finding biological important commonalities from a set of sequences. Many versions have been proposed and a huge number of papers have been written on effective and efficient methods for constructing multiple sequence alignment. We will discuss some of the important versions such as *SP-alignment*, *star alignment*, *tree alignment*, *generalized tree alignment*, and *fixed topology alignment with recombination*. Recent results on those versions are given.

We assume that the reader has the basic knowledge of algorithms and computational complexity (such as NP, P and MAX-SNP). Consult, *e.g.*, [27, 38, 57] otherwise.

The rest of this chapter is organized as follows. In Section 2, we discuss construction and comparison methods for evolutionary trees. In Section 3, we discuss briefly various distances for comparing sequences and explain in details the syntenic distance measure. In Section 4, we discuss multiple sequence alignment problems. We conclude in Section 5 with a few open problems.

2 Construction and Comparison of Evolutionary Trees

The evolution history of organisms is often conveniently represented as trees, called *phylogenetic trees* or simply *phylogenies*. Such a tree has uniquely labeled leaves and unlabeled interior nodes, can be *unrooted* or *rooted* if the evolutionary origin is known, and usually has internal nodes of degree 3. Figure 1 shows an example of a phylogeny. A phylogeny may also have *weights* on its edges, where an edge weight (more popularly known as *branch length* in genetics) could represent the evolutionary distance along the edge. Many phylogeny reconstruction methods, including the distance and maximum likelihood methods, actually produce weighted phylogenies. Figure 1 also shows a weighted phylogeny (the weights are for illustrative purposes only).

2.1 Phylogenetic Construction Methods

Phylogenetic construction methods use the knowledge of evolution of molecules to infer the evolutionary history of the species. The knowledge of evolution

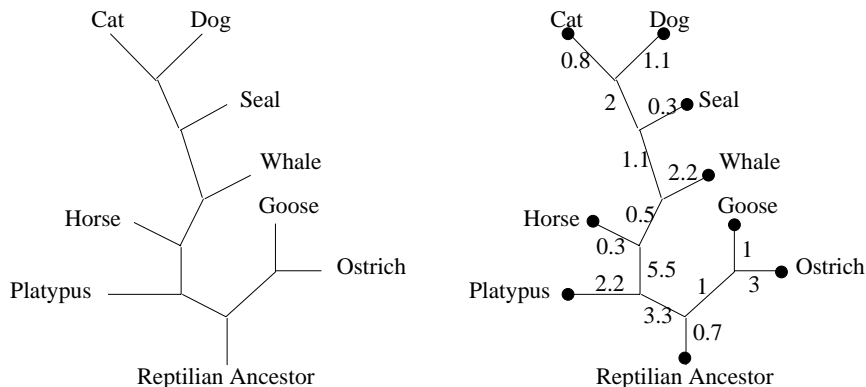


Figure 1: Examples of unweighted and weighted phylogenies.

is usually in the form of two kinds of data commonly used in phylogeny inference – namely, character matrices (where each position (i, j) is base j in sequence i), and distance matrices (where each position (i, j) contains the computed distance between sequence i and sequence j). Three major types of phylogenetic construction methods are the *parsimony and compatibility method*, the *distance method* and the *maximum-likelihood method*. Below we discuss each of them very briefly. See the excellent survey in [20, 72] for more details.

Parsimony methods construct phylogenetic trees for the given sequences such that, in some sense, the total number of changes (i.e., base substitutions) or some weighted sum of the changes is minimized. See [18, 22, 62] for some of the papers in this direction.

Distance methods [10, 23, 61] try to fit a tree to a matrix of pairwise distances between a set of n species. Entries in the distance matrices are assumed to represent evolutionary distance between species represented by the sequences in the tree, i.e., the total number of mutations in both lineages since divergence from the common ancestor. If no tree fits the distance matrix perfectly, then a measure of the discrepancy of the distances in the distance matrix and those in the tree is taken, and the tree with the minimum discrepancy is selected as the best tree. An example of the measure of the discrepancy, which has been used in the literature [10, 23], is a weighted least-square measure, i.e., of the form

$$\sum_{1 \leq i, j \leq n} w_{ij} (D_{ij} - d_{ij})^2$$

where D_{ij} are the given distances and d_{ij} are the distances computed from the tree.

Maximum-likelihood methods [18, 19, 8] relies on the statistical method of choosing a tree that maximizes the likelihood, i.e., maximizes the probability that the observed data would have occurred. Although this method is quite general and powerful, it is computationally intensive because of the complexity of the likelihood function.

All the above methods have been investigated by simulation and theoretical analysis. None of the methods work well under all evolutionary conditions, but each works well under particular situations. Hence, one must choose the appropriate phylogeny construction method carefully for best results [37].

2.2 Comparing Evolutionary Trees

As discussed in the previous section, over the past few decades, many approaches for reconstructing evolutionary trees have been developed, including (not exhaustively) parsimony, compatibility, distance and maximum-likelihood methods. As a result, in practice they often lead to different trees on the same set of species [48]. It is thus of interest to compare evolutionary trees produced by different methods, or by the same method on different data. Several distance models for evolutionary trees have been proposed in the literature. Among them, the best known is perhaps the *nearest neighbor interchange* (nni) distance introduced independently in [60] and [56]. Other distances include the *subtree-transfer* distance introduced in [34, 35] and the *linear-cost subtree-transfer distance* [14, 15]. Below, we discuss very briefly a few of these distances.

2.3 Nearest Neighbor Interchange Distance

An *nni* operation swaps two subtrees that are separated by an internal edge (u, v) , as shown in Figure 2. The *nni* operation is said to *operate* on this internal edge. The *nni* distance, $D_{nni}(T_1, T_2)$, between two trees T_1 and T_2 is defined as the minimum number of *nni* operations required to transform one tree into the other. K. Culik II and D. Wood [13] (improved later by [51]) proved that $n \log n + O(n)$ *nni* moves are sufficient to transform a tree of n leaves to any other tree with the same set of leaves. D. Sleator, R. Tarjan, and W. Thurston [69] proved an $\Omega(n \log n)$ lower bound for most pair of trees. Although the distance has been studied extensively in the literature

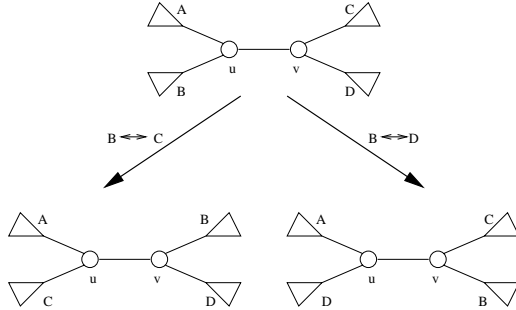


Figure 2: The two possible nni operations on an internal edge (u, v) : exchange $B \leftrightarrow C$ or $B \leftrightarrow D$.

[60, 56, 83, 13, 17, 40, 41, 43, 69, 51], the computational complexity of computing it has puzzled the research community for nearly 25 years until recently when the authors in [14] showed this problem to be NP-hard (an erroneous proof of the NP-hardness of the nni distance between unlabeled trees was published in [43]). Since computing the nni distance is shown to be NP-hard, the next obvious question is: *can we get a good approximation of the distance?* The authors in [51] show that the nni distance can be approximated in polynomial time within a factor of $\log n + O(1)$.

2.4 Subtree-transfer Distances

An nni operation can also be viewed as moving a subtree past a neighboring internal node. A more general operation is to transfer a subtree from one place to another arbitrary place. Figure 3 shows such a *subtree-transfer* operation. The subtree-transfer distance, $D_{st}(T_1, T_2)$, between two trees T_1

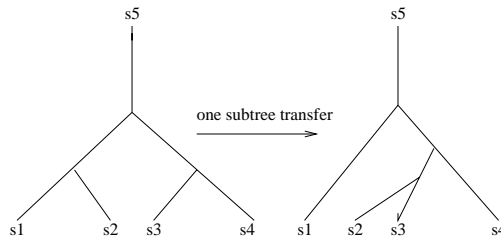


Figure 3: An example of a subtree-transfer operation on a tree.

and T_2 is the minimum number of subtrees we need to move to transform

T_1 into T_2 [34, 35, 36, 14, 15].

It is sometimes appropriate in practice to discriminate among subtree-transfer operations as they occur with different frequencies. In this case, we can charge each subtree-transfer operation a cost equal to the distance (the number of nodes passed) that the subtree has moved in the current tree. The *linear-cost* subtree-transfer distance, $D_{lcost}(T_1, T_2)$, between two trees T_1 and T_2 is then the minimum total cost required to transform T_1 into T_2 by subtree-transfer operations [14, 15]. Clearly, both subtree-transfer and linear-cost subtree-transfer models can also be used as alternative measures for comparing evolutionary trees generated by different tree reconstruction methods. In fact, on unweighted phylogenies, the linear-cost subtree-transfer distance is identical to the nni distance [15].

The authors in [36] show that computing the subtree-transfer distance between two evolutionary trees is NP-hard and give an approximation algorithm for this distance with performance ratio 3.

2.5 Rotation Distance

Rotation distance is a variant of the nni distance for rooted, ordered trees. A *rotation* is an operation that changes one rooted binary tree into another with the same size. Figure 4 shows the general rotation rule. An easy approximation algorithm for computing distance with a performance ratio of 2 is given in [68]. However, it is not known if computing this distance is NP-hard or not.

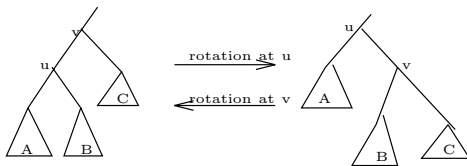


Figure 4: Left and right rotation operations on a rooted binary tree.

2.6 Distances on Weighted Phylogenies

Comparison of weighted evolutionary trees has recently been studied in [48]. The distance measure adopted is based on the difference in the partitions of the leaves induced by the edges in both trees, and has the drawback of being

somewhat insensitive to the tree topologies. Both the linear-cost subtree-transfer and nni models can be naturally extended to weighted trees. The extension for nni is straightforward: an nni is simply charged a cost equal to the weight of the edge it operates on. In the case of linear-cost subtree-transfer, although the idea is immediate, *i.e.* a moving subtree should be charged for the weighted distance it travels, the formal definition needs some care and can be found in [15].

Since computing the nni distance on unweighted phylogenies is NP-hard, it is obvious that computing this distance is NP-hard for weighted phylogenies also. The authors in [15] give an approximation algorithm for the linear-cost subtree-transfer distance on weighted phylogenies with performance ratio 2. In [14], the authors give an approximation algorithm for the nni distance on weighted phylogenies with performance ratio of $O(\log n)$. It is open whether the linear-cost subtree-transfer problem is NP-hard for weighted phylogenies. However, it has been shown that the problem is NP-hard for weighted trees with non-uniquely labeled leaves [15].

3 Computing Distances Between Genomes

The definition and study of appropriate measures of distance between pairs of species is of great importance in computational biology. Such measures of distance can be used, for example, in phylogeny construction and in taxonomic analysis.

As more and more molecular data becomes available methods for defining distances between species have focused on such data. One of the most popular distance measures is the edit distance between homologous DNA or aminoacid sequences obtained from different species. Such measures focus on point mutations and define the distance between two sequences as the minimum number of these moves required to transform one sequence into another. It has been recognized that the edit-distance may underestimate the distance between two sequences because of the possibility that multiple point mutations occurring at the same locus will be accounted for simply as one mutation. The problem is that the probability of a point mutation is not low enough to rule out this possibility.

Recently, there has been a spate of new definitions of distance that try to treat rarer, macrolevel mutations as the basic moves. For example, if we know the order of genes on a chromosome for two different species, we can define the *reversal* distance between the two species to be the number of

reversals of portions of the chromosome to transform the gene order in one species to the gene order in the other species. The question of finding the reversal distance was first explored in the computer science context by Kececioglu and Sankoff and by Bafna and Pevzner and there has been significant progress made on this question by Bafna, Hannenhalli, Kececioglu, Pevzner, Ravi, Sankoff and others [5, 6, 31, 45, 46]. Other moves besides reversals have been considered as well. Breaking off a portion of the chromosome and inserting it elsewhere in the chromosome is referred to as a *transposition* and one can similarly define the transposition distance[7]. Similarly allowing two chromosomes (viewed as strings of genes) to exchange suffixes (or sometimes a suffix with a prefix) is known as a *translocation* and this move can also be used to define an appropriate measure of distance between two species for which much of the genome has been mapped [44].

Ferretti et. al.[21] proposed a distance measure that is at an even higher level of abstraction. Here even the order of genes on a particular chromosome of a species is ignored/ presumed to be unknown. It is assumed that the genome of a species is given as a collection of sets. Each set in the collection corresponds to a set of genes that are on one chromosome and different sets in the collection correspond to different chromosomes (see Figure 5). In this

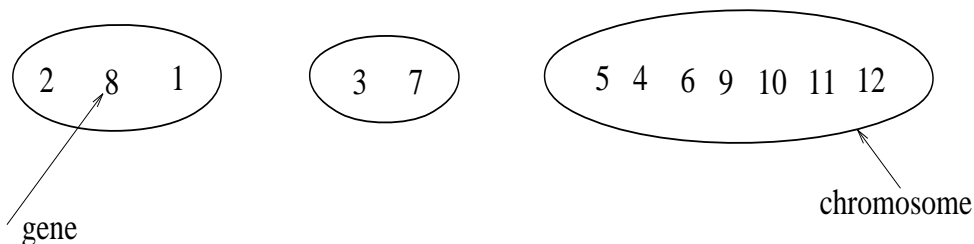


Figure 5: A Genome with 12 genes and 3 chromosomes

scenario one can define a move to be either an exchange of genes between two chromosomes, the fission of one chromosome into two, or the fusion of two chromosomes into one (see Figure 6). The *syntenic distance* between two species has been defined by Ferretti et. al.[21] to be the number of such moves required to transform the genome of one species to the genome of the other.

Notice that any recombination of two chromosomes is permissible in this model. By contrast, the set of legal translocations (in the translocation distance model) is severely limited by the order of genes on the chromosomes being translocated. Furthermore, the transformation of the first genome

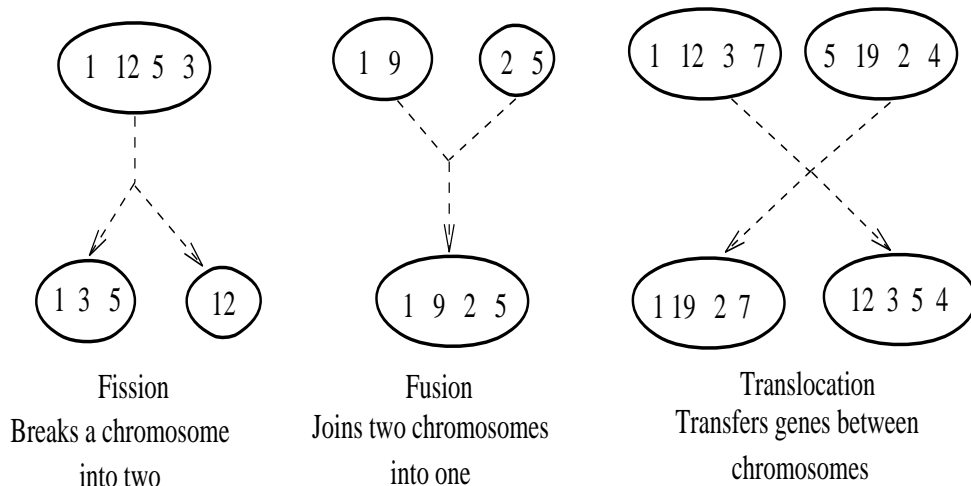


Figure 6: Different mutation operations

into the second genome does not have to produce a specified order of genes in the second genome. The underlying justification of this model is that the exchange of genes between chromosomes is a much rarer event than the movement of genes within a chromosome and hence a distance function should measure the minimum number of such exchanges needed.

In [16], the authors prove various results on the syntenic distance. For example, they show that computing the syntenic distance exactly is NP-hard, there is a simple polynomial time approximation algorithm for the synteny problem with performance ratio 2 and computing the syntenic distance is fixed parameter tractable.

The median problem arises in connection with the phylogenetic inference problem [21] and defined as follows. Given three genomes \mathcal{G}_1 , \mathcal{G}_2 and \mathcal{G}_3 , we are required to construct a genome \mathcal{G} such that the *median distance* $\alpha_{\mathcal{G}} = \sum_{i=1}^3 D(\mathcal{G}, \mathcal{G}_i)$ is minimized (where D is the syntenic distance). Without any additional constraints, this problem is trivial, since we can take \mathcal{G} to be empty (and then $\alpha_{\mathcal{G}} = 0$). In the context of syntenic distance, any one of the following three constraints seem relevant: (c1) \mathcal{G} must contain all genes present in *all the three* given genomes, (c2) \mathcal{G} must contain all genes present in *at least two* of the three given genomes, (c3) \mathcal{G} must contain all genes present in *at least one* of the three given genomes. Then, computing the median genome is NP-hard with any one of the three constraints **(c1)**, **(c2)** or **(c3)**. Moreover, one can approximate the median problem in polynomial

time (under any one of the constraints (c1), (c2) or (c3)) with a constant performance ratio. See [16] for details.

4 Multiple Sequence Alignment Problems

Multiple sequence alignment is the most critical cutting-edge tool for sequence analysis. It can help extracting, finding and representing biologically important commonalities from a set of sequences. These commonalities could represent some highly conserved subregions, common functions, or common structures. Multiple sequence alignment is also very useful in inferring the evolutionary history of a family of sequences [11, 29, 64, 82].

A *multiple alignment* \mathcal{A} of $k \geq 2$ sequences is obtained as follows: spaces are inserted into each sequence so that the resulting sequences s'_i ($i = 1, 2, \dots, k$) have the same length l , and the sequences are arranged in k rows of l columns each.

The value of the multiple alignment \mathcal{A} is defined as

$$\sum_{i=1}^l \mu(s'_1(i), s'_2(i), \dots, s'_k(i)),$$

where $s'_i(i)$ denotes the i -th letter in the resulting sequence s'_i , and $\mu(s'_1(i), s'_2(i), \dots, s'_k(i))$ denotes the score of the i -th column. The multiple sequence alignment problem is to construct a multiple alignment minimizing its value.

Many versions have been proposed based on different objective functions. We will discuss some of the important ones.

4.1 SP-alignment and Steiner consensus string

For *SP-score* (Sum-of-the-Pairs), the score of each column is defined as:

$$\mu(s'_1(i), s'_2(i), \dots, s'_k(i)) = \sum_{1 \leq j < l \leq k} \mu(s'_j(i), s'_l(i)),$$

where $\mu(s'_j(i), s'_l(i))$ is the score of the two opposing letters $s'_j(i)$ and $s'_l(i)$. The SP-score is sensible and has previously been studied extensively.

SP-alignment problem is to find an alignment with the smallest SP-score. It is first studied in [9] and subsequently used in [1, 3, 30, 58]. SP-alignment problem can be solved exactly by using dynamic programming. However, if there are k sequences and the length of sequences is n , it takes $O(n^k)$ time. Thus, it works for only small numbers of sequences. Some techniques

to reduce the time and space have been developed in [1, 28, 52, 66]. With these techniques, it is possible to optimally align up to 6 sequences of 200 characters in practice.

In fact, SP-alignment problem is NP-hard [74]. Thus, it is impossible to have a polynomial time algorithm for this problem. In the proof of NP-hardness, it is assumed that some pairs of identical characters have non-zero score. An interesting open problem is what if each pair of two identical characters is scored 0.

The first approximation algorithm was given by Gusfield [30]. He introduced the *center star* algorithm. Center star algorithm is very simple and efficient. It selects a sequence (called *center string*) s_c in the set of k given sequences S such that $\sum_{i=1}^k \text{dist}(s_c, s_i)$ is minimized. It then optimally aligns the sequences in $S - \{s_c\}$ to s_c and gets $k - 1$ pairwise alignments. These $k - 1$ pairwise alignments lead to a multiple alignment for the k sequences in S . If the score scheme for pairs of characters satisfies the triangle inequality, the cost of the multiple alignment produced by the center star algorithm is at most twice of the optimum [30, 29]. Some improved results were reported in [4, 58].

Another score called *consensus* score is defined as follows:

$$\mu(s'_1(i), s'_2(i), \dots, s'_k(i)) = \min_{s \in \Sigma} \sum_{j=1}^k \mu(s'_j(i), s),$$

where Σ is the set of characters that form the sequences. Here we reconstruct a character for each column and thus obtain a string. This string is called a *Steiner consensus string* and can be used as a representative for the set of given sequences. The problem is called the *Steiner consensus string* problem.

The Steiner consensus string problem was proved to be NP-complete [73] and MAX SNP-hard [74]. In the proof of MAX SNP-hardness, it is assumed that there is a “wild card”, and thus the triangle inequality does not hold. Combining with the results in [2], it shows that there is no polynomial time approximation scheme for this problem. Interestingly, the same center star algorithm also has performance ratio 2 for this problem [29].

4.2 Tree alignment

Tree score: In order to define the score $\mu(s'_1(i), s'_2(i), \dots, s'_k(i))$ of the i -th column, an *evolutionary* (or *phylogenetic*) tree $T = (V, E)$ with k leaves is assumed, each leaf j corresponding to a sequence s_j . (Here V and E denote

the sets of nodes and edges in T , respectively.) Let $k + 1, k + 2, \dots, k + m$ be the internal nodes of T . For each internal node j , reconstruct a letter (possibly a space) $s'_j(i)$ such that $\sum_{(p,q) \in E} \mu(s'_p(i), s'_q(i))$ is minimized. The score $\mu(s'_1(i), s'_2(i), \dots, s'_k(i))$ of the i -th column is thus defined as

$$\mu(s'_1(i), s'_2(i), \dots, s'_k(i)) = \sum_{(p,q) \in E} \mu(s'_p(i), s'_q(i)).$$

This measure has been discussed in [1, 4, 62, 63, 64]. Multiple sequence alignment with tree score is often referred to as *tree alignment* in the literature.

Note that, a tree alignment induces a set of *reconstructed* sequences, each corresponding to an internal node. Thus, it is convenient to reformulate tree alignment as follows: Given a set X of k sequences and an evolutionary tree T with k leaves, where each leaf is associated with a given sequence, reconstruct a sequence for each internal node to minimize the *cost* of T . Here, the cost of T is the sum of the edit distance of each pair of (given or reconstructed) sequences associated with an edge. Observe that, once a sequence for each internal node has been reconstructed, a multiple alignment can be obtained by optimally aligning the pair of sequences associated with each edge of the tree. Moreover, the tree score of this induced multiple alignment equals the cost of T . In this sense, the two formulations of tree alignment are equivalent.

Sankoff gave an exact algorithm for tree alignment that runs in $O(n^k)$, where n is the length of the sequences and k is the number of given sequences. Tree alignment was proved to be NP-hard [74].

Therefore it is unlikely to have a polynomial time algorithm for tree alignment. Some heuristic algorithms have also been considered in the past. Altschul and Lipman tried to cut down the computation volume required by dynamic programming [1]. Sankoff, Cedergren and Lapalme gave an iterative improvement method to speed up the computation [63, 64]. Waterman and Perlwitz devised a heuristic method when the sequences are related by a binary tree [80]. Hein proposed a heuristic method based on the concept of a *sequence graph* [32, 33]. Ravi and Kecericioglu designed an approximation algorithm with performance ratio $\frac{deg+1}{deg-1}$ when the given tree is a *regular deg-ary* tree (i.e., each internal node has exactly deg children) [59].

The first approximation algorithm with a guaranteed performance ratio was devised by Wang, Jiang, and Lawler [75]. A ratio-2 algorithm was given. The algorithm was then extended to a polynomial time approximation

scheme (PTAS), i.e., the performance ratio could arbitrarily approach 1. The PTAS requires computing exact solutions for depth- t subtrees. For a fixed t , the performance ratio was proved to be $1 + \frac{3}{t}$, and the running time was proved to be $O((k/deg^t)^{deg^{t-1}+2}M(2, t-1, n))$, where deg is the degree of the given tree, n is the length of the sequences, and $M(deg, t-1, n)$ is the time needed to optimally align a tree with $deg^{t-1} + 1$ leaves, which is upper-bounded by $O(n^{deg^{t-1}+1})$. Based on the analysis, to obtain a performance ratio less than 2, exact solutions for depth-4 subtrees must be computed, and thus optimally aligning 9 sequences at a time is required. This is impractical even for sequences of length 100.

An improved version was given in [76]. They proposed a new PTAS for the case where the given tree is a regular deg -ary tree. The algorithm is much faster than the one in [75]. The algorithm also must do local optimizations for depth- t subtrees. For a fixed t , the performance ratio of the new PTAS is $1 + \frac{2}{t} - \frac{2}{t2^t}$ and the running time is $O(\min\{2^t, k\}kdM(deg, t-1, n))$, where d is the depth of the tree. Presently, there are efficient programs [63] to do local optimizations for three sequences ($t = 2$). In fact, we can expect to obtain optimal solutions for 5 sequences ($t = 3$) of length 200 in practice since there is such a program [28, 52] for SP-score and similar techniques can be used to attack tree alignment problem. Therefore, solutions with costs at most 1.583 times the optimum can be obtained in practice for strings of length 200.

For tree alignment, the given tree is typically a binary tree. Recently, Wang, Jiang and Gusfield design a PTAS for binary trees. The new approximation scheme adopts a more clever partitioning strategy and has a better time efficiency for the same performance ratio. For any fixed r , where $r = 2^{t-1} + 1 - q$ and $0 \leq q \leq 2^{t-2} - 1$, the new PTAS runs in time $O(kdn^r)$ and achieves an approximation ratio of $\frac{2^{t-1}}{2^{t-2}(t+1)-q}$. Here the parameter r represents the “size” of local optimization. In particular, when $r = 2^{t-1} + 1$, its approximation ratio is simply $\frac{2}{t+1}$.

4.3 Generalized Tree alignment

In practice, we often face a more difficult problem called *generalized tree alignment*. Suppose we are given a set of sequences. The problem is to construct an evolutionary tree as well as a set of sequences (called reconstructed sequences) such that each leaf of the evolutionary tree is assigned a given sequence, each internal node of the tree is assigned a reconstructed

sequence, and the cost of the tree is minimized over all possible evolutionary trees and reconstructed sequences.

Intuitively, the problem is harder than tree alignment since the tree is not given and we have to compute the tree structure as well as the sequences assigned to internal nodes. In fact, the problem was proved to be MAX SNP-hard [74] and a simplified proof was given in [78]. It implies that it is impossible to have a PTAS for generalized tree alignment unless P=NP [2]. This confirms the observation from approximation point of view.

Generalized tree alignment problem is in fact the Steiner tree problem in sequence spaces. One might use the approximation algorithms with guaranteed performance ratios [84] for graph Steiner trees. However, this may lead to a tree structure where a given sequence is an internal node. Sometimes, it is unacceptable. Schwikowski and Vingron give a method that combines clustering algorithms and Hein’s sequence graph method. The produced solutions contain biologically reasonable trees and keep the guaranteed performance ratio. [67].

4.4 Fixed topology history/alignment with recombination

Multigene families, viruses, and alleles from within populations experience recombinations [34, 35, 47, 71]. When recombination happens, the ancestral material on the present sequence s_1 is located on two sequences s_2 and s_3 . s_2 and s_3 can be cut at k locations (break points) into $k + 1$ pieces, where $s_2 = s_{2,1}s_{2,2} \dots s_{2,l+1}$ and $s_3 = s_{3,1}s_{3,2} \dots s_{3,l+1}$. s_1 can be represented as $s_{2,1}\hat{s}_{3,2}s_{2,3} \dots s_{2,i}\hat{s}_{3,i+1} \dots$, where subsequences $s_{2,i}$ and $s_{3,i+1}$ differ from the corresponding $s_{2,i}$ and $s_{3,i+1}$ by insertion, deletion, and substitution of letters. k , the number of times s_1 switches between s_2 and s_3 , is called the number of *crossovers*. The cost of the recombination is

$$\begin{aligned} dist(s_{1,1}, s_{1,1}) + dist(s_{2,2}, s_{2,2}), \dots dist(s_{1,i}, s_{1,i}) + dist(s_{2,i+1}, s_{2,i+1}) \\ + \dots + k\chi \end{aligned}$$

where $dist(s_{2,i+1}, s_{2,i+1})$ is the edit distance between the two sequences $s_{2,i+1}$ and $s_{2,i+1}$, k is the number of crossovers and χ is the crossover penalty. The *recombination* distance to produce s_1 from s_2 and s_3 is the cost of a recombination that has the smallest cost among all possible recombinations. We use $r_dist(s_1, s_2, s_3)$ to denote the recombination distance. For more details, see [47, 81].

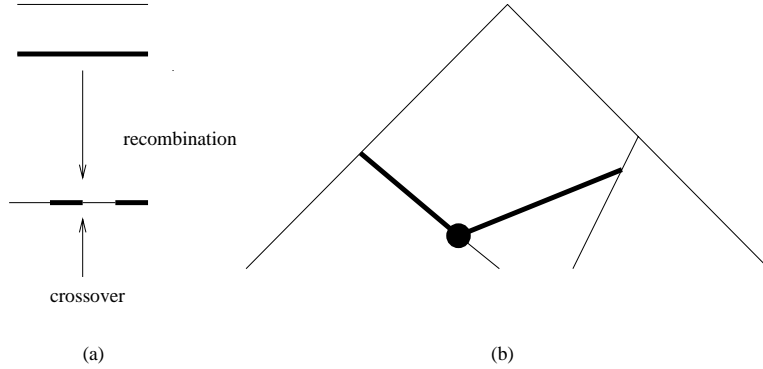


Figure 7: (a) Recombination operation. (b) The topology. The dark edges are recombination edges. The circled node is a recombination node.

When recombination occurs, the given topology is no longer a binary tree. Instead, some nodes, called *recombination nodes*, in the given topology may have two parents[34, 35]. In a more general case as described in [47], the topology may have more than one root. The set of roots is called a *pro-toset*. The edges incident to recombination nodes are called *recombination edges*. See Figure 7 (b). A node/edge is *normal* if it is not a recombination node/edge.

The cost of a pair of recombination edges is the recombination distance to produce the sequence on the recombination node from the two sequences on its parents. The cost of other normal edges is the edit distance between two sequences. A topology is *fully labeled*, if every node in the topology is labeled. For a fully labeled topology, the cost of the topology is the total cost of edges in the topology. Each node in the topology with degree greater than 1 is an internal node. Each leaf/terminal (degree 1 node) in the topology is labeled with a given sequence. The goal here is to construct a sequence for each internal node such that the cost of the topology is minimized. We call this problem *fixed topology history with recombination* (FTHB).

Obviously, this problem is a generalization of tree alignment. The difference is that the given topology is no longer a binary tree. Instead, there are some recombination nodes which have two parents instead of one. Moreover, there may be more than one root in the topology.

A different version called *fixed topology alignment with recombination*

(FTAR) is also discussed [53]. From approximation point of view, FTHR and FTAR are much harder than tree alignment. It is shown that FTHR and FTAR cannot be approximated within any constant performance ratio unless $P = NP$ [53].

A more restricted case, where each internal node has at most one recombination child and there are at most 6 parents of recombination nodes in any path from the root to a leaf in the given topology, is also considered. It is shown that the restricted version for both FTHR and FTAR is MAX-SNP-hard. That is, there is no polynomial time approximation scheme unless $P = NP$ [53].

The above hardness results are disappointing. However, recombination occur infrequently. So, it is interesting to study some restricted cases. A *merge node* of recombination node v is the lowest common ancestor of v 's two parents. The two different paths from a recombination node to its merge node are called *merge paths*. We then study the case, where

- (C1) each internal node has at most one recombination child and
- (C2) any two merge paths for different recombination nodes do not share any common node.

Using a method similar to the lifting method for tree alignment, one can get a ratio-3 approximation algorithm for both FTHR and HTAR when the given topology satisfies (C1) and (C2). The ratio-3 algorithm can be extended to a PTAS for FTAR with bounded number of crossovers. (See [53].)

Remarks: Hein might be the first to study the method to reconstruct the history of sequences subject to recombination [34, 35]. Hein observed that the evolution of a sequence with k recombinations could be described by k recombination points and $k + 1$ trees describing the evolution of the $k + 1$ intervals, where two neighboring trees were either identical or differed by one subtree transfer operation [34, 35, 36, 15, 14]. A heuristic method was proposed to find the most parsimonious history of the sequences in terms of mutation and recombination operations.

Another strike was given by Kececioğlu and Gusfield [47]. They introduced two new problems, *recombination distance*, and *bottleneck recombination history*. They tried to include higher-order evolutionary events such as block insertions and deletions [26], and tandem repeats [42, 49].

5 Conclusion

In this chapter we have discussed some important topics in the field of computational biology such as the phylogenetic construction and comparison methods, syntenic distance between genomes and the multiple sequence alignment problems. Given the vast majority of topics in computational biology, these discussed topics constitute only a part of them. Some of the important topics which were *not* covered in this chapter are:

- protein structure prediction,
- DNA physical mapping problems,
- metabolic modeling,
- string / database search problems etc.

We hope that this survey article will inspire the readers for further study and research of these and other related topics.

Papers on computational molecular biology have started to appear in many different books, journals and conferences. Below we list some sources which could serve as excellent starting points for various problems that arise in computational biology:

Books: References [12, 30, 39, 54, 65, 70, 82].

Journals: Computer Applications in the Biosciences (recently renamed as Bioinformatics), Journal of Computational Biology, Bulletin of Mathematical Biology, Journal of Theoretical Biology etc.

Conferences: Annual Symposium on Combinatorial Pattern Matching (CPM), Pacific Symposium on Biocomputing (PSB), Annual International Conference on Computational Molecular Biology (RECOMB), Annual Conference on Intelligent Systems in Molecular Biology (ISMB) etc.

Web pages: <http://www.cs.washington.edu/education/courses/590bi>,
<http://www.cse.ucsc.edu/research/compbio>,
<http://www.cs.jhu.edu/~salzberg/cs439.html> etc.

Acknowledgments

We thank Prof. Tao Jiang for bringing the authors together. Thanks also go to Dr. Todd Wareham who carefully reads the draft and gives valuable suggestions.

References

- [1] S. Altschul and D. Lipman. Trees, stars, and multiple sequence alignment, *SIAM Journal on Applied Math.*, 49 (1989), pp. 197-209.
- [2] S. Arora, C. Lund, R. Motwani, M. Sudan and M. Szegedy. On the intractability of approximation problems, *33rd IEEE Symposium on Foundations of Computer Science*, 1992, pp. 14-23.
- [3] D. Baconn and W. Anderson. Multiple sequence alignment, *Journal of Moleccular Biology*, 191 (1986), pp. 153-161.
- [4] V. Bafna, E. Lawer and P. Pevzner. Approximate methods for multiple sequence alignment, *Proc. 5th Symp. on Combinatorial Pattern Matching. Springer LNCS 807*, 1994, pp. 43-53.
- [5] V. Bafna and P. Pevzner. Genome Rearrangements and Sorting by Reversals, *34th IEEE Symp. on Foundations of Computer Science*, 1993, pp. 148-157.
- [6] V. Bafna and P. Pevzner. Sorting by Reversals: Genome Rearrangements in Plant Organelles and Evolutionary History of X Chromosome, *Mol. Biol. and Evol.*, 12 (1995), pp. 239-246.
- [7] V. Bafna and P. Pevzner. Sorting by Transpositions, *Proc. of 6th Ann. ACM-SIAM Symp. on Discrete Algorithms*, 1995, pp. 614-623.
- [8] D. Barry and J.A. Hartigan. Statistical analysis of hominoid molecular evolution, *Stat. Sci.*, 2 (1987), pp. 191-210.
- [9] H. Carrillo and D. Lipman. The multiple sequence alignment problem in biology, *SIAM Journal on Applied Mathematics*, 48 (1988), pp. 1073-1082.
- [10] L. L. Cavalli-Sforza and A. W. F. Edwards. Phylogenetic analysis: models and estimation procedures, *Evolution*, 32 (1967), pp. 550-570. (also published in *Am. J. Hum. Genet.*, 19, pp. 233-257)
- [11] S. C. Chan, A. K. C. Wong and D. K. T. Chiu. A survey of multiple sequence comparison methods, *Bulletin of Mathematical Biology*, 54, 4 (1992), pp. 563-598.

- [12] J. Collado-Vides, B. Magasanik and T. F. Smith (eds.). *Integrative Approaches to Molecular Biology*, MIT Press; Cambridge, MA, 1996.
- [13] K. Culik II and D. Wood. A note on some tree similarity measures, *Information Processing Letters*, 15 (1982), pp. 39-42.
- [14] B. DasGupta, X. He, T. Jiang, M. Li, J. Tromp and L. Zhang. On distances between phylogenetic trees, *Proc. 8th Annual ACM-SIAM Symposium on Discrete Algorithms*, 1997, pp. 427-436.
- [15] B. DasGupta, X. He, T. Jiang, M. Li, and J. Tromp. On the linear-cost subtree-transfer distance, to appear in the special issue in *Algorithmica* on computational biology, 1998.
- [16] B. DasGupta, T. Jiang, S. Kannan, M. Li and E. Sweedyk. On the Complexity and Approximation of Syntenic Distance, *1st Annual International Conference On Computational Molecular Biology*, 1997, pp. 99-108 (journal version to appear in *Discrete and Applied Mathematics*).
- [17] W. H. E. Day. Properties of the nearest neighbor interchange metric for trees of small size, *Journal of Theoretical Biology*, 101 (1983), pp. 275-288.
- [18] A.W.F. Edwards and L.L. Cavalli-Sforza. The reconstruction of evolution, *Ann. Hum. Genet.*, 27 (1964), 105. (Also in *Heredity* 18, 553.)
- [19] J. Felsenstein. Evolutionary trees for DNA sequences: a maximum likelihood approach, *J. Mol. Evol.*, 17 (1981), pp. 368-376.
- [20] J. Felsenstein. Phylogenies from molecular sequences: inferences and reliability, *Annu. Rev. Genet.*, 22 (1988), pp. 521-565.
- [21] V. Ferretti, J.H. Nadeau and D. Sankoff. Original Synteny. In *Proc. of 7th Ann. Symp. on Combinatorial Pattern Matching*, 1996, pp. 159-167.
- [22] W.M. Fitch. Toward defining the course of evolution: minimum change for a specified tree topology, *Syst. Zool.*, 20 (1971), pp. 406-416.
- [23] W.M. Fitch and E. Margoliash. Construction of phylogenetic trees, *Science*, 155 (1967), pp. 279-284.
- [24] K. A. Frenkel. The human genome project and informatics, *Communications of the ACM*, 34, 11 (1991), pp. 41-51.

- [25] V. A. Funk and D. R. Brooks. *Phylogenetic Systematics as the Basis of Comparative Biology*, Smithsonian Institution Press; Washington, DC, 1990.
- [26] Z. Galil and R. Ciancarlo. Speeding up dynamic programming with applications to molecular biology, *Theoretical Computer Science*, 64 (1989), pp. 107-118.
- [27] M. R. Garey and D. S. Johnson. *Computers and Intractability: A Guide to the Theory of NP-completeness*, W. H. Freeman, 1979.
- [28] S. Gupta, J. Kececioğlu, and A. Schaffer. Making the shortest-paths approach to sum-of-pairs multiple sequence alignment more space efficient in practice, *Proceedings of the 6th Symposium on Combinatorial Pattern Matching, Springer LNCS937*, 1995, pp. 128-143.
- [29] D. Gusfield. *Algorithms on Strings, Trees, and Sequences: Computer Science and Computational Biology*, Cambridge University Press, 1997.
- [30] D. Gusfield. Efficient methods for multiple sequence alignment with guaranteed error bounds, *Bulletin of Mathematical Biology*, 55 (1993), pp. 141-154.
- [31] S. Hannenhalli and P. Pevzner. Transforming Cabbage into Turnip (polynomial algorithm for sorting signed permutations by reversals), *Proc. of 27th Ann. ACM Symp. on Theory of Computing*, 1995, pp. 178-189.
- [32] J. Hein. A tree reconstruction method that is economical in the number of pairwise comparisons used, *Mol. Biol. Evol.*, 6, 6 (1989), pp. 669-684.
- [33] J. Hein. A new method that simultaneously aligns and reconstructs ancestral sequences for any number of homologous sequences, when the phylogeny is given, *Mol. Biol. Evol.*, 6 (1989), pp. 649-668.
- [34] J. Hein. Reconstructing evolution of sequences subject to recombination using parsimony, *Math. Biosci.*, 98 (1990), pp. 185-200.
- [35] J. Hein. A heuristic method to reconstruct the history of sequences subject to recombination, *J. Mol. Evol.*, 36 (1993), pp. 396-405.
- [36] J. Hein, T. Jiang, L. Wang, and K. Zhang. On the complexity of comparing evolutionary trees, *Discrete Applied Mathematics*, 71(1996), 153-169.

- [37] D. M. Hillis, B. K. Mable and C. Moritz. Applications of Molecular Systematics, in D. M. Hillis, C. Moritz and B. K. Mable (eds.), *Molecular Systematics (Second Edition)*, Sinauer Associates, Sunderland, MA, 1996, pp. 515-543.
- [38] D. Hochbaum. *Approximation Algorithms for NP-hard Problems*, PWS publishers, 1996.
- [39] L. Hunter (ed.), *Artificial Intelligence in Molecular Biology*, MIT Press, Cambridge MA, 1993.
- [40] J. P. Jarvis, J. K. Luedeman and D. R. Shier. Counterexamples in measuring the distance between binary trees, *Mathematical Social Sciences*, 4 (1983), pp. 271-274.
- [41] J. P. Jarvis, J. K. Luedeman and D. R. Shier. Comments on computing the similarity of binary trees, *Journal of Theoretical Biology*, 100 (1983), pp. 427-433.
- [42] S. Kannan and E. W. Myers. An algorithm for locating non-overlapping regions of maximum alignment score, *3rd Annual Symposium on Combinatorial Pattern Matching*, 1993, pp. 74-86.
- [43] M. Křivánek. Computing the nearest neighbor interchange metric for unlabeled binary trees is NP-complete, *Journal of Classification*, 3 (1986), pp. 55-60.
- [44] J. Kececioglu and R. Ravi. Of Mice and Men: Evolutionary Distances Between Genomes under Translocation, *Proc. of 6th Ann. ACM-SIAM Symp. on Discrete Algorithms*, 1995, pp. 604-613.
- [45] J. Kececioglu and D. Sankoff. Exact and Approximation Algorithms for the Inversion Distance between Two Permutations, *Proc. of 4th Ann. Symp. on Combinatorial Pattern Matching*, Lecture Notes in Computer Science 684, Springer Verlag, 1993, pp. 87-105.
- [46] J. Kececioglu and D. Sankoff. Efficient Bounds for Oriented Chromosome Inversion Distance, *Proc. of 5th Ann. Symp. on Combinatorial Pattern Matching*, Lecture Notes in Computer Science 807, Springer Verlag, 1994, pp. 307-325.

- [47] J. Kececioglu and D. Gusfield. Reconstructing a history of recombinations from a set of sequences, 5th Annual ACM-SIAM Symposium on Discrete Algorithms, 1994, pp. 471-480.
- [48] M. Kuhner and J. Felsenstein. A simulation comparison of phylogeny algorithms under equal and unequal evolutionary rates. *Mol. Biol. Evol.* 11, 3 (1994), pp. 459-468.
- [49] G. M. Landau and J. P. Schmidt. An algorithm for approximate tandem repeats, 3rd Ann. Symp. on Combinatorial Pattern Matching, 1993, pp. 120-133.
- [50] E.S. Lander, R. Langridge and D.M. Saccocio. Mapping and interpreting biological information, *Communications of the ACM*, 34, 11 (1991), pp. 33-39.
- [51] M. Li, J. Tromp, and L.X. Zhang. On the nearest neighbor interchange distance between evolutionary trees, *Journal of Theoretical Biology*, 182 (1996), pp. 463-467.
- [52] J.Lipman, S.F. Altschul, and J.D. Kececioglu. A tool for multiple sequence alignment, *Proc. Nat. Acad. Sci. U.S.A.*, 86, pp.4412-4415, 1989.
- [53] B. Ma, L. Wang and M. Li. Fixed topology alignment with recombination, *CPM98*, to appear.
- [54] J. Meidanis and J. C. Setubal. *Introduction to Computational Molecular Biology*, PWS Publishing Company, Boston, MA, 1997.
- [55] W.Miller, S. Schwartz, and R. C. Hardison. A point of contact between computer science and molecular biology, *IEEE Computational Science and Engineering*, Spring 1994, pp. 69-78.
- [56] G. W. Moore, M. Goodman and J. Barnabas. An iterative approach from the standpoint of the additive hypothesis to the dendrogram problem posed by molecular data sets, *Journal of Theoretical Biology*, 38 (1973), pp. 423-457.
- [57] C. H. Papadimitriou. *Computational Complexity*, Addison-Wesley; reading, MA, 1994.
- [58] P. Pevzner. Multiple alignment, communication cost, and graph matching, *SIAM Journal of Applied Mathematics*, 56, 6 (1992), pp. 1763-1779.

- [59] R. Ravi and J. Kececioglu. Approximation algorithms for multiple sequence alignment under a fixed evolutionary tree, *5th Annual Symposium on Combinatorial Pattern Matching*, 1995, pp. 330-339.
- [60] D. F. Robinson. Comparison of labeled trees with valency three, *Journal of Combinatorial Theory, Series B*, 11 (1971), pp. 105-119.
- [61] N. Saitou and M. Nei. The neighbor-joining method: a new method for reconstructing phylogenetic trees, *Mol. Biol. Evol.*, 4 (1987), pp. 406-425.
- [62] D. Sankoff. Minimal mutation trees of sequences, *SIAM Journal of Applied Mathematics*, 28 (1975), pp. 35-42.
- [63] D. Sankoff, R. J. Cedergren and G. Lapalme. Frequency of insertion-deletion, transversion, and transition in the evolution of 5S ribosomal RNA, *J. Mol. Evol.* 7 (1976), pp. 133-149.
- [64] D. Sankoff and R. Cedergren. Simultaneous comparisons of three or more sequences related by a tree, in D. Sankoff and J. Kruskal, editors, *Time warps, string edits, and macromolecules: the theory and practice of sequence comparison*, pp. 253-264, Addison Wesley, 1983.
- [65] D. Sankoff and J. Kruskal (eds.). *Time warps, string edits, and macromolecules: the theory and practice of sequence comparison*, Addison Wesley, 1983.
- [66] G.D. Schuler, S.F. Altschul, and D.J. Lipman. A workbench for multiple alignment construction and analysis, in *Proteins: Structure, function and Genetics*, in press.
- [67] B. Schwikowski and M. Vingron. The deferred path heuristic for the generalized tree alignment problem, *1st Annual International Conference On Computational Molecular Biology*, 1997, pp. 257-266.
- [68] D. Sleator, R. Tarjan, W. Thurston. Rotation distance, triangulations, and hyperbolic geometry, *J. Amer. Math. Soc.*, 1 (1988), pp. 647-681.
- [69] D. Sleator, R. Tarjan, W. Thurston. Short encodings of evolving structures, *SIAM J. Discr. Math.*, 5 (1992), pp. 428-450.
- [70] G. A. Stephens. *String Searching Algorithms*, World Scientific Publishers, Singapore, 1994.

- [71] F. W. Stahl. *Genetic recombination*, Scientific American, pp. 90-101, February 1987.
- [72] D. L. Swofford, G. J. Olsen, P. J. Waddell and D. M. Hillis. *Phylogenetic Inference*, in D. M. Hillis, C. Moritz and B. K. Mable (eds.), *Molecular Systematics* (Second Edition), Sinauer Associates, Sunderland, MA, 1996, pp. 407-514.
- [73] E. Sweedyk and T. Warnow, The tree alignment problem is NP-complete, *Manuscript*.
- [74] L. Wang and T. Jiang. On the complexity of multiple sequence alignment, *Journal of Computational Biology*, 1 (1994), pp. 337-348.
- [75] L. Wang, T. Jiang and E.L. Lawler. Approximation algorithms for tree alignment with a given phylogeny, *Algorithmica*, 16 (1996), pp. 302-315.
- [76] L. Wang and D. Gusfield. Improved approximation algorithms for tree alignment, *Journal of Algorithms*, 25 (1997), pp. 255-173.
- [77] L. Wang, T. Jiang, and Dan Gusfield. A more efficient approximation scheme for tree alignment, *1st Annual International Conference on Computational Molecular Biology*, 1997, pp. 310-319.
- [78] H. T. Wareham, A simplified proof of the NP-hardness and MAX SNP-hardness of multiple sequence tree alignment, *Journal of Computational Biology*, Vol. 2, pp. 509-514, 1995.
- [79] M.S. Waterman. Sequence alignments, in *Mathematical Methods for DNA Sequences*, M.S. Waterman (ed.), CRC, Boca Raton, FL, 1989, pp. 53-92.
- [80] M.S. Waterman and M.D. Perlwitz. Line geometries for sequence comparisons, *Bull. Math. Biol.*, 46 (1984), pp. 567-577.
- [81] J. D. Watson, N. H. Hopkins, J. W. Roberts, J. A. Steitz, A. M. Weiner. *Molecular Biology of the gene*, 4th edition, Benjamin-Cummings, Menlo Park, California, 1987.
- [82] M.S. Waterman. *Introduction to Computational Biology: Maps, sequences, and genomes*, Chapman and Hall, 1995.
- [83] M. S. Waterman and T. F. Smith. On the similarity of dendrograms, *Journal of Theoretical Biology*, 73 (1978), pp. 789-800.

- [84] A.Z. Zelikovsky. The $11/6$ approximation algorithm for the Steiner problem on networks, *Algorithmica*, 9 (1993), pp. 463-470.