

## 5.1 Approximation Algorithms for Multiple Sequence Alignment

### 5.1.1 Problem Definition

Given a family  $\mathcal{S} = (S_1, \dots, S_k)$  of  $k$  sequences, such that the sequences are “similar” to each other, we would like to find out the common characteristics of this family. Aligning each pair of sequences from  $\mathcal{S}$  separately, often does not reveal this common information.

A *multiple alignment* of  $\mathcal{S}$  is a new set of sequences  $\mathcal{S}' = (S'_1, \dots, S'_k)$  such that:

- All the strings in  $\mathcal{S}'$  are of equal length. We denote this length by  $l$ .
- Each  $S'_i$  was generated from  $S_i$  by inserting spaces.

When performing multiple alignment, as in the case of pairwise alignment, one wishes to evaluate the quality of the alignment by giving it a numeric score (see also lecture 3).

**Definition 5.1** - *The sum of pairs (SP) score of a multiple alignment  $\mathcal{M}$  is the sum of the scores of pairwise global alignments induced by  $\mathcal{M}$ .*

Let  $\sigma(x, y)$  be our scoring function, i.e., the price of aligning the character  $x$  with the character  $y$ , for  $x, y \in \Sigma \cup \{-\}$ . We assume that  $\sigma(-, -) = 0$ ,  $\sigma(x, y) = \sigma(y, x)$ , and that the triangle inequality  $\sigma(x, y) \leq \sigma(x, z) + \sigma(z, y)$  holds.

### 5.1.2 The Center Star Method for (SP) Alignment

In this section, we present an approximation algorithm for calculating the optimal multiple alignment under the SP metric (see e.g. [3] [pp 348-350]). The algorithm achieves an approximation ratio of two.

Given a multiple alignment  $\mathcal{M}$ , let  $d(S_i, S_j)$  be the score of the pairwise alignment it induces on  $S_i, S_j$ . Our target function is the *SP value* of  $\mathcal{M}$  which is  $d(\mathcal{M}) \equiv \sum_{i < j} d(S_i, S_j)$  (sum over all pairs of the score of the induced alignment).

---

<sup>1</sup>Based in part on a scribe by Sarel Har-Peled on December 18, 1995 and on the book [3, pp 347-363]

**Problem 5.1** *The SP alignment problem.*

**INPUT:** *A set of sequences  $S$ .*

**QUESTION:** *Compute a global multiple alignment  $\mathcal{M}$  with minimum sum-of-pairs score.*

We denote by  $D(S, Y)$  the score of the optimal alignment between sequences  $S$  and  $Y$ .

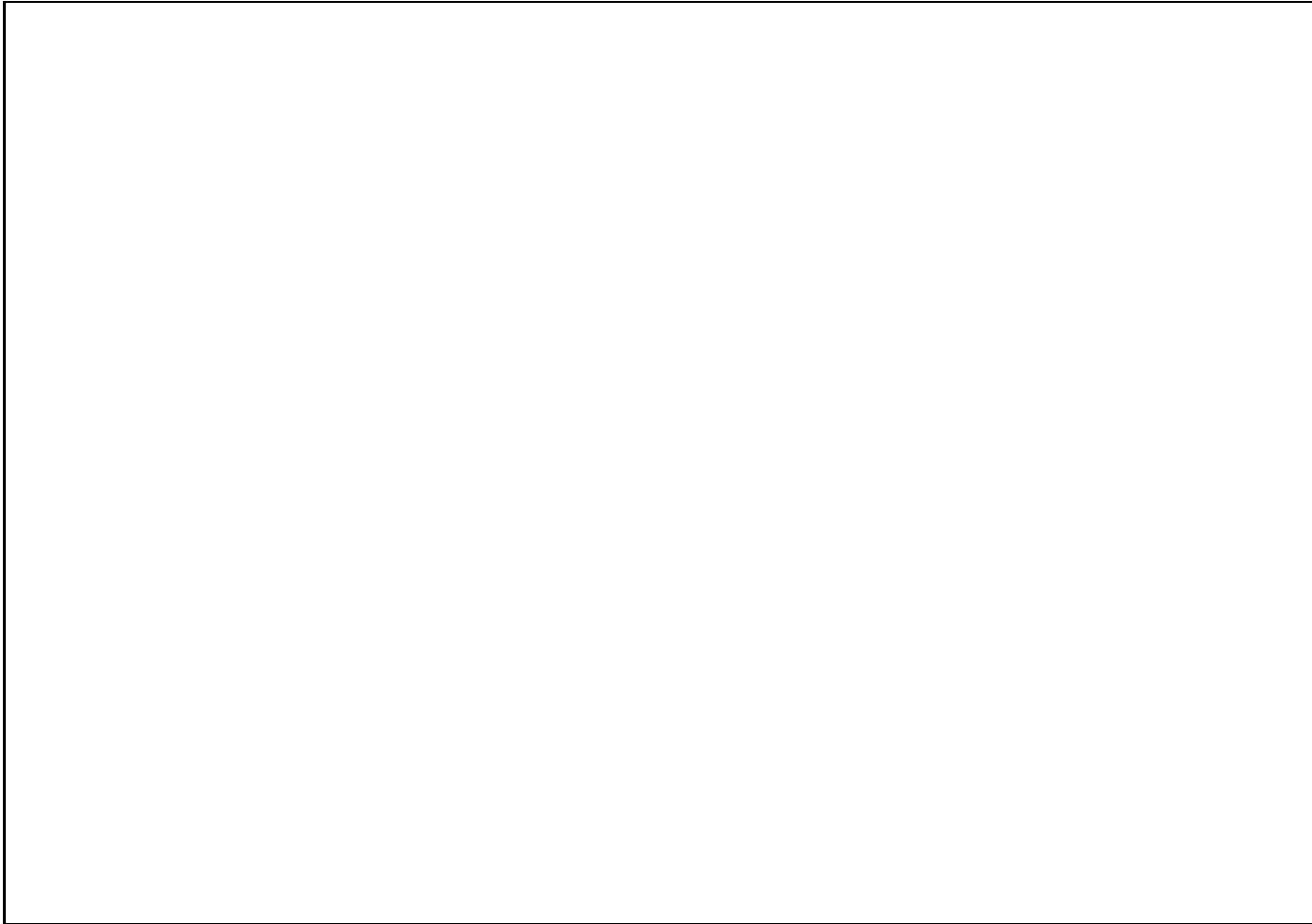


Figure 5.1: A generic center star for six strings, where the center string ( $S_c$ ) is  $S_3$ .

**Definition 5.2** *We say that a tree  $T$  having the elements of  $S$  as its nodes, induces a multiple alignment  $\mathcal{M}(T)$  over  $S$ , if for each edge  $(S, Y) \in E(T)$  the value of the alignment induced by  $\mathcal{M}$  on  $S$  and  $Y$  is  $D(S, Y)$  (optimal). In such a case we write  $\mathcal{M} = \mathcal{M}(T)$ .*

Given  $S = (S_1, \dots, S_k)$ , let  $S_c \in S$  denote the element of  $S$ , for which  $\sum_{i \neq c} D(S_i, S_c)$  is minimal. We refer to  $S_c$  as the *center* of  $\mathcal{M}_c$ .

**Definition 5.3** We define the center star,  $T_c$  to be a star tree of  $k$  nodes, with the center node labeled  $S_c$  and with each of the  $k - 1$  remaining nodes labeled by a distinct sequence in  $\mathcal{S} \setminus \{S_c\}$  (see figure 5.1). The multiple alignment  $\mathcal{M}_c$  of  $\mathcal{S}$  is the multiple alignment induced by the center star, i.e., for each  $v \neq c$ , the alignment  $\mathcal{M}_c$  induces an optimal pairwise alignment between  $S_c$  and  $S_v$ .

**The Center Star Algorithm:**

1. Find  $S_t \in \mathcal{S}$  minimizing  $\sum_{i \neq t} D(S_i, S_t)$  and let  $\mathcal{M} = \{S_t\}$ .
2. Add the sequences in  $\mathcal{S} \setminus \{S_t\}$  to  $\mathcal{M}$  one by one so that the alignment of every newly added sequence with  $S_t$  is optimal. Add spaces, when needed, to all pre-aligned sequences.

**Running time analysis:**

1.  $\binom{k}{2} O(n^2)$  for step 1.
2.  $\sum_{i=1}^{k-1} O((i \cdot n) \cdot n) = O(k^2 \cdot n^2)$  for step 2. (Since the worst-case length of  $S'_c$  after the addition of  $i$  strings is  $(i + 1) \cdot n$ )

**Lemma 5.2** For  $1 \leq i, j \leq k, i \neq j$  it holds that  $d(S_i, S_j) \leq D(S_i, S_c) + D(S_c, S_j)$ .

*Proof:* Since the triangle inequality holds for every single column of the alignment by the definition of the scoring scheme, it also holds for entire strings by the definition of  $d$ . Therefore  $d(S_i, S_j) \leq d(S_i, S_c) + d(S_c, S_j)$ . But the edges  $(S_i, S_c), (S_c, S_j)$  are edges of  $E(T_c)$ , thus  $d(S_i, S_c) = D(S_i, S_c)$ . It follows that  $d(S_i, S_j) \leq D(S_i, S_c) + D(S_c, S_j)$ . ■

Let  $\mathcal{M}^*$  denote the optimal alignment of  $\mathcal{S}$ . Let  $d^*(S_i, S_j)$  denote the value of alignment between  $S_i$  and  $S_j$  induced by  $\mathcal{M}^*$ .

**Theorem 5.3**

$$\frac{d(\mathcal{M}_c)}{d(\mathcal{M}^*)} \leq \frac{2(k-1)}{k} < 2$$

*Proof:* By Lemma 5.2 it follows that

$$2d(\mathcal{M}_c) = \sum_{i \neq j} d(S_i, S_j) \leq \sum_{i \neq j} (D(S_i, S_c) + D(S_c, S_j)) = 2(k-1) \sum_{i \neq c} D(S_c, S_i) \quad (5.1)$$

We define  $W$  to be  $\sum_{i \neq c} D(S_c, S_i)$  and we get:

$$2d(\mathcal{M}_c) \leq 2(k-1)W \quad (5.2)$$

On the other hand, by the choice of  $c$  it follows that:

$$2d(\mathcal{M}^*) = \sum_{i \neq j} d^*(S_i, S_j) \geq \sum_{i \neq j} D(S_i, S_j) = \sum_i \sum_{j \neq i} D(S_i, S_j) \geq \sum_i W = kW. \quad (5.3)$$

And finally,

$$\frac{d(\mathcal{M}_c)}{d(\mathcal{M}^*)} \leq \frac{2(k-1)W}{kW} = \frac{2(k-1)}{k}. \quad (5.4)$$

■

Theorem 5.3 implies that calculating the multiple alignment of the center star produces a multiple alignment with a value which is at most  $R_k = \frac{2(k-1)}{k}$  times the value of the optimal alignment. For example  $R_3 = \frac{4}{3}$ ,  $R_4 = \frac{3}{2}$ .

### 5.1.3 Multiple Alignment with Consensus

In this section we look at an approximation algorithm for a multiple alignment that optimizes a different score metrics - the consensus error. As before, we assume the existence of a pairwise scoring scheme  $\sigma$  satisfying the triangle inequality.

**Definition 5.4** *Given a set of strings  $\mathcal{S}$ , the consensus error of a string  $\bar{S}$  with respect to  $\mathcal{S}$  is  $E(\bar{S}) = \sum_{S_i \in \mathcal{S}} D(\bar{S}, S_i)$ . Note that  $\bar{S}$  need not be in  $\mathcal{S}$ .*

**Problem 5.4** *Optimal Steiner string.*

**INPUT:** *A set of strings  $\mathcal{S}$*

**QUESTION:** *Find a string  $S^*$  which minimizes the consensus error  $E(S^*)$  over all possible strings.*

The Steiner string  $S^*$  attempts to capture the common characteristics of the set of strings  $\mathcal{S}$  and reflect them in a single string. We will present an approximation algorithm for the optimal Steiner string problem with worst-case approximation ratio of 2.

**Lemma 5.5** *Let  $\mathcal{S}$  have  $k$  strings, and assume that the scoring scheme  $\sigma$  satisfies the triangle inequality. Then there exists a string  $\bar{S} \in \mathcal{S}$  such that  $\frac{E(\bar{S})}{E(S^*)} \leq 2 - \frac{2}{k} < 2$  (see e.g. [3] [pp 349–351]).*

*Proof:* For any  $\bar{S} \in \mathcal{S}$ :

$$\begin{aligned} E(\bar{S}) &= \sum_{S_i \in \mathcal{S}} D(\bar{S}, S_i) \leq \sum_{S_i \neq \bar{S}} [D(\bar{S}, S^*) + D(S^*, S_i)] = & (5.5) \\ & (k-2) \cdot D(\bar{S}, S^*) + D(\bar{S}, S^*) + \sum_{S_i \neq \bar{S}} D(S^*, S_i) = (k-2) \cdot D(\bar{S}, S^*) + E(S^*) \end{aligned}$$

If we pick  $\bar{S} \in \mathcal{S}$  such that  $\bar{S}$  is closest to  $S^*$  then:

$$E(S^*) = \sum_{S_i \in \mathcal{S}} D(S^*, S_i) \geq k \cdot D(\bar{S}, S^*) \quad (5.6)$$

The *center string*  $S_c \in \mathcal{S}$  minimizes  $\sum_{S_i \in \mathcal{S}} D(S_c, S_i)$  and therefore its consensus error is smaller than the consensus error of the  $\bar{S}$  (the string closest to  $S^*$ ). We get:

$$\frac{E(S_c)}{E(S^*)} \leq \frac{(k-2) \cdot D(S_c, S^*) + E(S^*)}{E(S^*)} \leq \frac{(k-2) \cdot D(S_c, S^*)}{k \cdot D(S_c, S^*)} + 1 = 2 - \frac{2}{k} \quad (5.7)$$

The proof above uses the lemma 5.5 and the fact that  $E(S_c) \leq E(\bar{S})$  ■

It is worthwhile noting that Steiner string was defined without alignment, and the only requirement is the distance function, that satisfies the triangle inequality. We will next start discussing consensus strings that are alignment motivated.

### 5.1.4 Consensus Strings from Multiple Alignment

**Definition 5.5** *Given a multiple alignment  $\mathcal{M}$  of a set of strings  $\mathcal{S}$ , the consensus character in column  $i$  of  $\mathcal{M}$  is the character that minimizes the summed distance to it from all the characters in column  $i$ . Let  $d(i)$  denote that minimum sum in column  $i$ .*

**Definition 5.6** *The consensus string  $S_{\mathcal{M}}$  derived from the alignment  $\mathcal{M}$  is the concatenation of the consensus characters for each column of  $\mathcal{M}$ .*

**Definition 5.7** *The alignment error of  $S_{\mathcal{M}}$  equals  $\sum_{i=1}^l d(i)$  where  $l$  is the number of characters in  $S_{\mathcal{M}}$*

**Definition 5.8** *The optimal consensus multiple alignment is a multiple alignment  $\mathcal{M}$  of an input set  $\mathcal{S}$  whose consensus string  $S_{\mathcal{M}}$  minimizes the alignment error. It can be shown that the optimal consensus multiple alignment is equal to the optimal Steiner string, as defined in section 5.1.3.*

We can use the center string ( $S_c$ ) for approximating the optimal multiple alignment with an alignment error smaller than  $(2 - \frac{2}{k})$  times the optimal alignment error.

## 5.2 Multiple Alignment to a Phylogenetic Tree

**Definition 5.9** *A tree  $T$  with a distinct string label (from a set of string  $\mathcal{S}$ ) assigned to each leaf is called a phylogenetic tree on  $\mathcal{S}$*

**Definition 5.10** Given a phylogenetic tree  $T$  on  $\mathcal{S}$ , a phylogenetic alignment  $T'$  for  $T$  is an assignment of one string label to each internal node of  $T$ . Note that the strings assigned to internal nodes need not be distinct and need not be from the set  $\mathcal{S}$ .

For example, when  $T$  is a star, choosing the sequence to label its central node is a phylogenetic alignment.

The phylogenetic tree  $T$  is meant to represent the "established" evolutionary history of a set of objects of interest, with the convention that each extant object is represented at a unique leaf of the tree. Each edge  $(u, v)$  represents some evolutionary history that transforms the string at  $u$  (assuming  $u$  is the parent of  $v$ ) to the string at  $v$ . For convenience, when denoting an edge by a pair of nodes, we shall hereafter write the parent node first.

**Definition 5.11** If strings  $S$  and  $S'$  are assigned to the endpoints of an edge  $(i, j)$ , then the edge distance of  $(i, j)$  is defined to be  $D(S, S')$ .

**Definition 5.12** Let  $\mathcal{M}_T$  be a phylogenetic alignment for a phylogenetic tree  $T$ . The distance of  $\mathcal{M}_T$  is given by the sum of all the edge distances over all the edges of  $T$ .

**Problem 5.6** *Phylogenetic alignment:*

**INPUT:** A set  $\mathcal{S} = (S_1, \dots, S_k)$  of strings and a phylogenetic tree  $T$ , on  $\mathcal{S}$ .

**QUESTION:** find a phylogenetic alignment  $\mathcal{M}_T$  with minimum distance.

We assume that the structure of the tree  $T$  is known to us. This usually happens in practice when  $T$  was previously reconstructed using solid evolutionary data.

### 5.2.1 Lifted Alignment tree - a Heuristic for Phylogenetic Alignment

**Definition 5.13** A phylogenetic alignment is called lifted alignment if for every internal node  $v$ , the string assigned to  $v$  is also assigned to one of  $v$ 's children (See figure 5.2). Trivially, in such a case, all internal nodes are assigned labels from the set of leaf strings.

Let  $T^*$  be the optimal alignment for tree  $T$ . We will construct a lifted alignment  $T^L = \text{Lift}(T^*)$ , which is based on  $T^*$ , with only a limited damage to the alignment distance. Note that this construction is only conceptual since we usually do not know  $T^*$ .

For each node  $v$  let its  $T^*$  label be  $S_v^*$ . We shall assign every  $v$  a label  $S_v^L$ . Initially, only the leaves are labeled, and by definition  $S_v^L = S_v^*$  for each leaf  $v$ . The labeling process successively traverses the internal nodes in any order, provided that a node is not visited before any of its children. Upon visiting a node, it is *lifted* i.e. labeled by one of the labels of its children. Thus, the resulting phylogenetic alignment is lifted. (see figure 5.3)

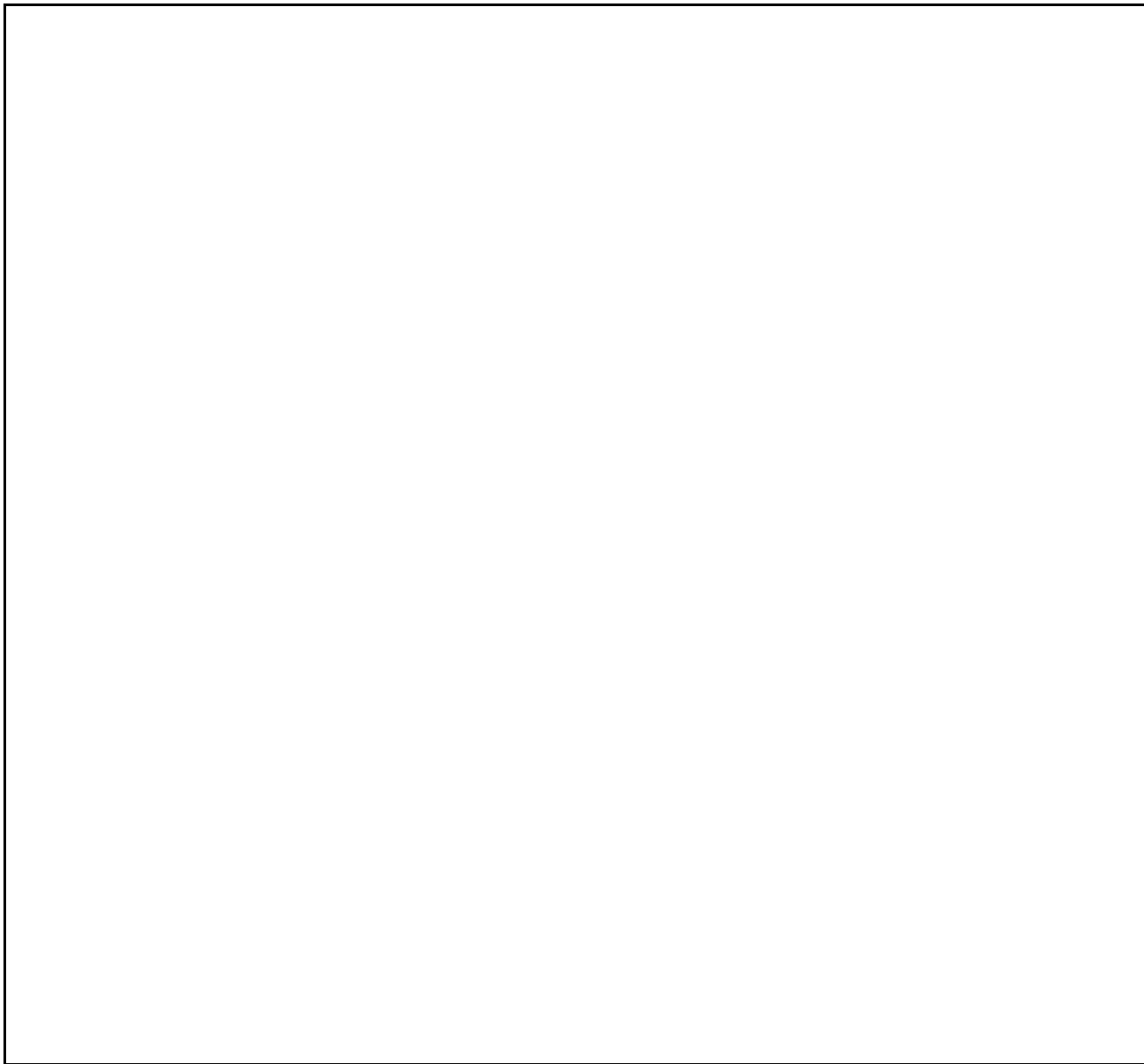


Figure 5.2: A phylogenetic tree with lifted alignment. Each internal node is labeled by one of the strings labeling its children.

```

Procedure Lift( $T$  : Tree)
  begin
    while there exists an unlifted node  $v$ , all of whose children have been lifted, do :
      Find a child  $j$  whose label  $S_j$  is the closest to  $S_v^*$ . Namely
        For every child  $i$  of  $v$ :
           $D(S_v^*, S_j) \leq D(S_v^*, S_i)$ 
        Label  $S_v^*$  with  $S_j$ 
    end while
  end

```

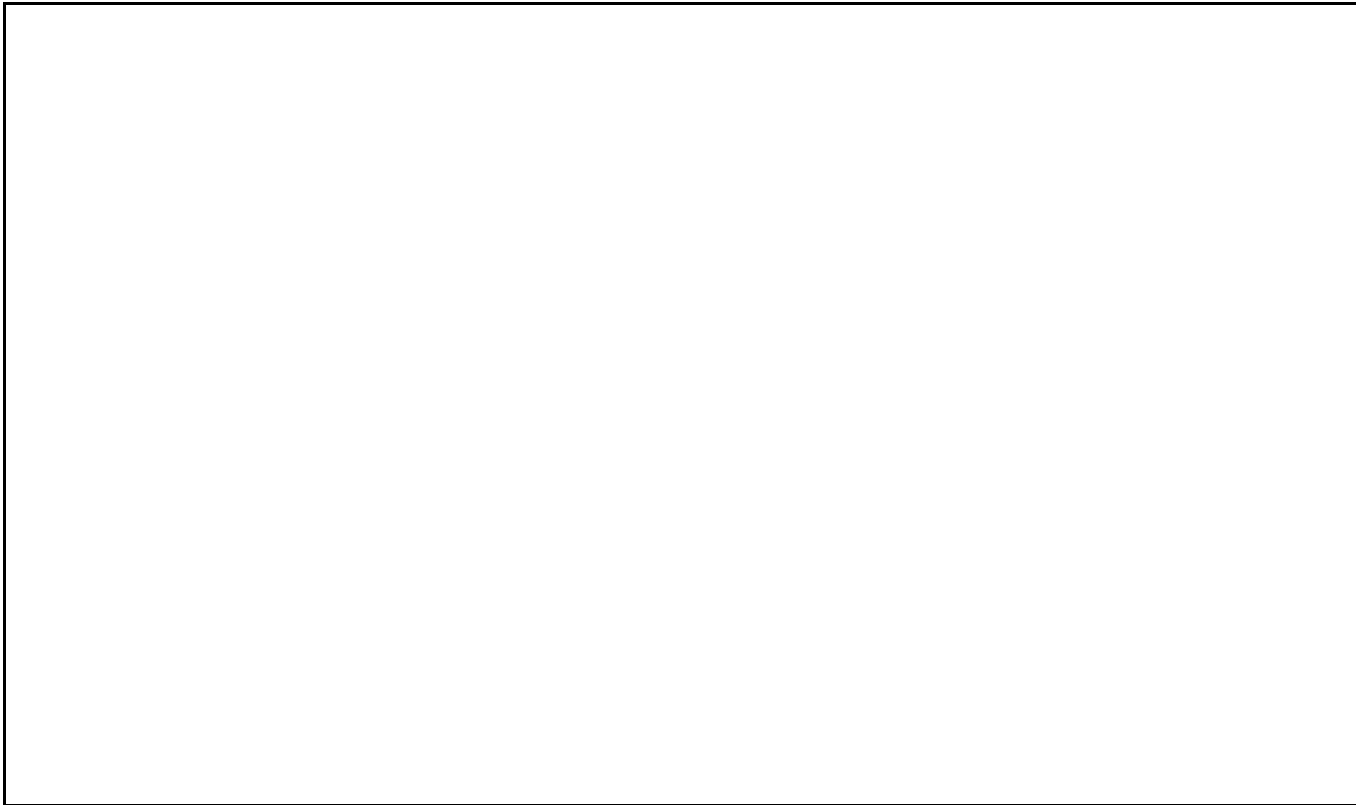


Figure 5.3: The lifting construction at node  $v$ . The numbers on the edges are the distances from  $S_v^*$  to the lifted strings labeling its children. On the left is the tree before lifting, and on the right the result of the lift. After the lift one edge will have a distance of 0.

**Theorem 5.7** (Jiang, Wang and Lawler, 1996 [5]) *The distance of the phylogenetic alignment  $T^L = Lift(T^*)$  is at most twice the distance of the optimal phylogenetic alignment  $T^*$ .*

*Proof:* Let  $e = (v, w)$  be an edge in  $T$ . Suppose that in  $T^L$ ,  $S_j$  is the label of  $v$  and  $S_i$  is the label of  $w$ . If  $i = j$  then  $D(S_j, S_i) = 0$ . Otherwise:

$$D(S_i, S_j) \leq D(S_j, S_v^*) + D(S_v^*, S_i) \leq 2 \cdot D(S_v^*, S_i) \quad (5.8)$$

The first inequality is due to the triangle inequality, and the second follows from the labeling algorithm. For an edge  $e = (v, w)$ , with  $S_w = S_i$ , Let  $P_e$  be the path in  $T$  from  $v$  to the leaf labeled  $S_i$ . Due to the triangle inequality

$$D(S_v^*, S_i) \leq \text{the total length of } P_e \text{ in } T^* \quad (5.9)$$

We say that the edge  $e = (v, w)$  is *blue* in  $T^L$  if  $S_i \neq S_j$ . The distance of a lifted alignment  $T^L$  is equal to the sum of edge distances on all the blue edges in the tree.

For a blue edge  $e = (v, w)$ , observe that the definition of lifted alignment implies that along the path  $P_e$  every node except  $v$  is labeled  $S_i$ , and no node outside  $P_e$  is labeled  $S_i$ . Hence, if  $e' = (v', w')$  is any other blue edge, then  $P_e$  and  $P_{e'}$  have no edges in common. This defines a mapping from every blue edge  $e$  in  $T^L$  to a path  $P_e$  in  $T^*$  such that:

- The distance in  $T^L$  of the edge  $e$  is at most twice the total distance in  $T^*$  of the edges on  $P_e$ . (Follows from equations 5.8 and 5.9.)
- No edge in  $T^*$  is mapped to by more than one edge in  $T^L$ .

Therefore the total distance of  $T^L$  = the total distance on blue edges  $\leq 2 \cdot$  the sum of all total distances of  $P_e$  paths in  $T^* \leq 2 \cdot$  the total distance of  $T^*$ . (see figure 5.4) ■

We now describe how to find the optimal lifted alignment using a dynamic programming algorithm as listed below. But first we define:

**Definition 5.14** Let  $T_v$  be the subtree of  $T$  rooted at node  $v$  and  $S \in \mathcal{S}$ . Let  $d(v, S)$  denote the distance of the best lifted alignment of  $T_v$  under the requirement that string  $S$  is assigned to node  $v$ .

The algorithm will compute  $d(v, S)$  for any  $S \in \mathcal{S}$  working it's way from the leaves up using the following recursion:

- If  $v$  is an internal node with all its children being leaves, then  $d(v, S) = \sum_{(v,w)} D(S, S_w)$  where  $S_w$  is the label of  $w$ .
- Else,  $d(v, S) = \sum_{(v,v')} \min_{S' \in \mathcal{S}} [D(S, S') + d(v', S')]$  where  $v'$  is a child of  $v$  and  $S'$  is a label of one of the leaves of  $T_{v'}$ .

**Time analysis:** We perform a preprocessing stage, in which we compute all the  $\binom{k}{2}$  pairwise distances between the  $k$  input strings. This takes  $O(N^2)$  time, where  $N$  is the total length of all the strings. The work at any internal node is  $O(k^2)$ , and the overall work of the algorithm is  $O(N^2 + k^3)$

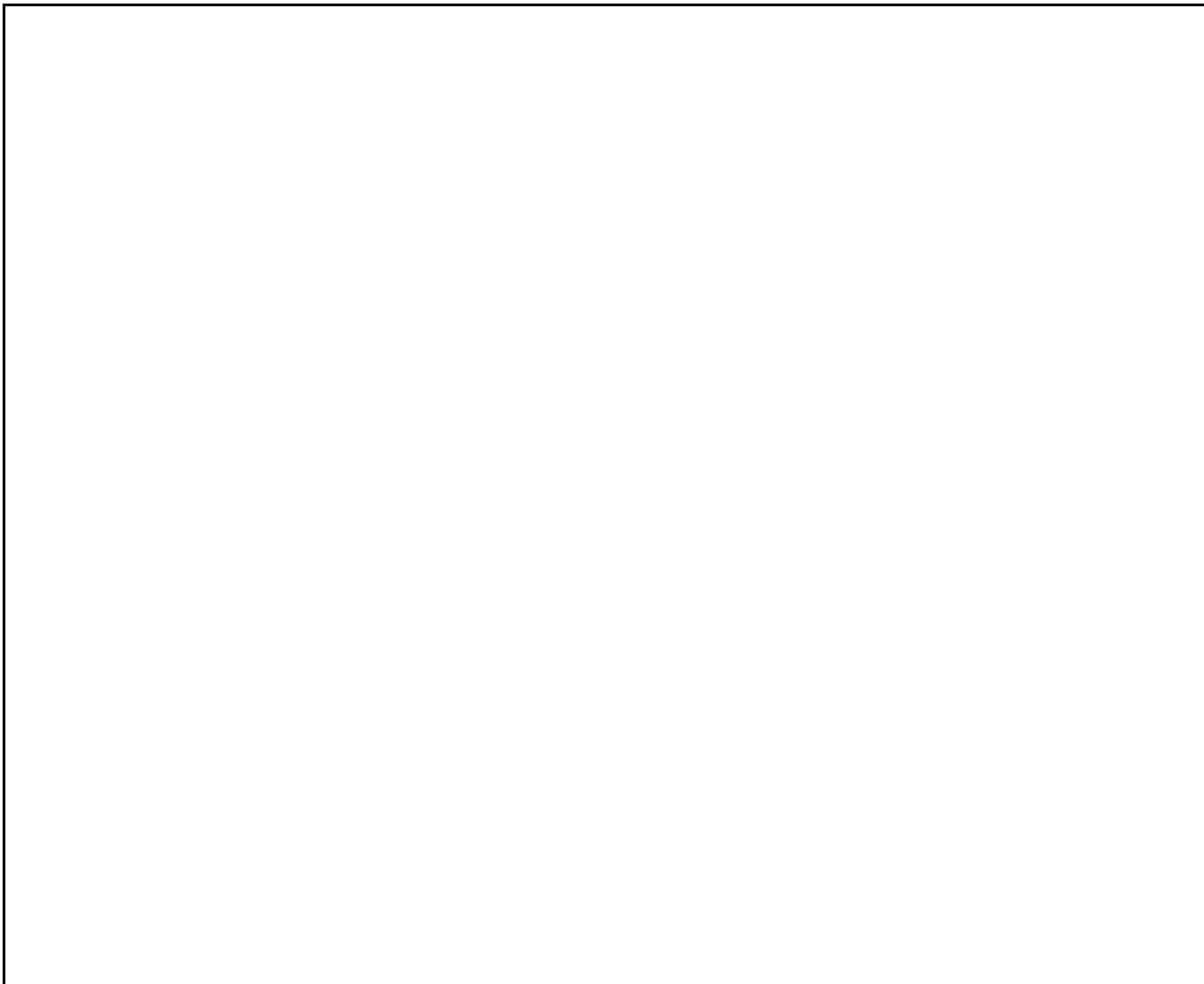


Figure 5.4: The lifted tree  $T_L^*$ . The dashed edges show the paths along which a leaf string has been lifted to some internal node. Solid edges are blue edges in  $T_L^*$ , while each of these dashed edges has distance 0. The path  $P_{(a,b)}$  for example, is the path  $b, d, S_4$  along which the string labeling  $b$  was lifted. Edge  $(a, b)$  has distance in  $T_L^*$  at most twice the distance of path  $P_{(a,b)}$  in  $T^*$ .

## 5.3 Common Multiple Alignments Methods

### 5.3.1 Aligning a String to a Profile

Given a database of sequences, we would like to partition it into families of “similar” sequences. For this purpose we would like to encompass our knowledge on the common properties of the sequences in a family *profile* (formal definition to follow). Constructing a profile of a family enables us to identify its members and test whether or not a new sequence belongs to the family. Moreover, searching the database with a profile is more sensitive than searching using a single sequence of the family: When searching with a single sequence, we can only look for the sequences in the database with the best alignments with the given sequence, while when using a profile we may test for membership to the family.

**Definition 5.15** *for an alignment  $S'$  of length  $l$ , a profile is a  $l \times |\Sigma \cup \{-\}|$  matrix, whose columns are probability vector denoting the frequencies of each symbol in the corresponding alignment column.*

Any alignment between a sequence  $B$  and a profile  $P$  (i.e. both have the same length) can be evaluated by  $\sum_{j=1}^m \sigma(p_j, b_j)$ . Clearly, using dynamic programming, we can find the best alignment of a sequence against a profile.

The key in pairwise alignment is scoring two positions  $x$  and  $y$  :  $\sigma(x, y)$ . For a letter  $x$  and a column  $y$  of a profile, let  $\sigma(x, y)$  be the probability of  $x$  being in column  $y$ . The value for  $x$  depends on the frequency of its occurrences in the column  $y$ . We also need to devise a score for  $\sigma(x, -)$ . In order to find whether a given sequence is a member of certain family, we use a usual pairwise dynamic programming alignment to compare the given sequence to the family profile.

### 5.3.2 Iterative pairwise alignment

This approach uses pairwise alignment scores to iteratively add one additional string to a growing multiple alignment. We start by aligning the two strings whose edit distance is the minimum over all pairs of strings. Then we iteratively consider the string with the smallest distance to any of the strings already in the multiple alignment.

More generally, the algorithm works as follows:

1. Align some pair.
2. while (not done)
  - (a) Pick an unaligned string which is "near" some aligned one(s).
  - (b) Align with the *profile* of the previously aligned group  
Resulting new spaces are inserted into all strings in the group.

### 5.3.3 Progressive alignment

#### Feng-Doolittle 1987

In this algorithm (Feng and Doolittle, 1987 [1]), the key idea is that the pair of strings with minimum distance is almost likely to have been obtained from the pair of objects that had most recently diverged, and that the pairwise alignment of these two specific strings provides the most "reliable" information that can be extracted from the input strings. Therefore any spaces (gaps) that appear in the optimal pairwise alignment of those two strings should be preserved in the overall multiple alignment.

The algorithm is as follows:

1. Calculate the  $\binom{k}{2}$  pairwise alignment scores, and convert them to distances.
2. Use an incremental clustering algorithm (Fitch and Margoliash, 1967 [2]) to construct a tree from the distances.
3. Traverse the nodes in their order of addition to the tree, repeatedly align the child nodes (sequences or alignments).

Features of this heuristic:

- Highest scoring pairwise alignment determines the alignment to two groups.
- "Once a gap, always a gap": replace gaps in alignments by a neutral character.

#### CLUSTALW

ClustalW is a software package for multiple alignment (implementing an algorithm of Thompson, Higgins, Gibson 1994 [4]). The basic idea is the same as in the Feng-Doolittle algorithm:

1. Calculate the  $\binom{k}{2}$  pairwise alignment scores, and convert to two distances.
2. Use clustering algorithm (Neighbor-Joining) to build a tree from the distances.
3. align sequence - sequence, sequence - profile, profile - profile in decreasing similarity order.

This algorithm makes use of many ad-hoc rules such as weighting, different matrix scores and special gap scores.

# Bibliography

- [1] D. Feng and R. F. Doolittle. Progressive sequence alignment as a prerequisite to correct phylogenetic trees. *J. Mol. Evol.*, 60:351–360, 1987.
- [2] W. M. Fitch and E. Margoliash. Construction of phylogenetic trees. *Science*, 155:279–284, 1967.
- [3] D. Gusfield. *Algorithms on Strings, Trees and Sequences*. Cambridge University Press, New York, 1997.
- [4] D. G. Higgins J. D. Thompson and T. J. Gibson. Clustal w: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice. *Nucleic Acids Res*, 22:4673–80, 1994.
- [5] E. L. Lawler T. Jiang, L. Wang. Approximation algorithms for tree alignment with a given phylogeny. *Algorithmica*, 16:302–315, 1996.