

# An Efficient Algorithm for Statistical Multiple Alignment on Arbitrary Phylogenetic Trees

G.A. Lunter, I. Miklós, Y.S. Song, J. Hein

December 12, 2002

*Department of Statistics, University of Oxford,  
1 South Parks Road, Oxford, OX1 3TG, UK*

## Abstract

We present an efficient algorithm for statistical multiple alignment based on the TKF91 model on an arbitrary  $k$ -leaved phylogenetic tree. The existing algorithms use a hidden Markov model approach, which requires at least  $O(\sqrt{5}^k)$  states and leads to a time complexity of  $O(5^k L^k)$ , where  $L$  is the geometric mean sequence length. Using a combinatorial technique reminiscent of inclusion/exclusion, we are able to sum away the states, thus improving the time complexity to  $O(2^k L^k)$  and considerably reducing memory requirements. This makes statistical multiple alignment under the TKF91 model a definite practical possibility in the case of a phylogenetic tree with a modest number of leaves.

## 1 Introduction

A major way by which biological sequences evolve is mutation. Three main types of mutation events are substitutions, insertions and deletions of amino acids or nucleotides. The latter two events, insertions and deletions, introduce the problem of *aligning* sequences, so that homologous positions appear in the same column of an alignment. When the aim is to find sequences from a database which are homologous to a query sequence, current alignment techniques perform quite well. The alignment problem in phylogenetics, however, is still a serious challenge [Lee01, Gol98]. For example, the phylogenies inferred from different, but equally good, alignments can be quite different [Gol98]. Some researchers have suggested that the regions in the sequences which are responsible for this variance in the inferred phylogeny (usually referred to as *unalignable regions* [Lee01]) should be omitted. It is unclear, however, how much information is lost by ignoring such regions. Moreover, several studies have suggested that accurate evolutionary parameters cannot be obtained using only a single “best” alignment [TKF91, HWK<sup>+</sup>00], even if this alignment is seemingly reliable [HWK<sup>+</sup>00]. A more robust approach is to take into account all possible alignments, or a large subset of those, in a statistical framework.

Since insertions and deletions, like substitutions, are rare random events, it seems natural to model them by a continuous time stochastic process, with rates for the three types of mutations. Stochastic models for the substitution process are well known and widely used to obtain maximum-likelihood evolutionary parameters [Fel81, WLG01]. The TKF91 [TKF91] model of sequence evolution incorporates such a mutation model, and moreover allows for nucleotides to spawn new nucleotides and to die. To model evolution on a phylogenetic tree, the TKF91 model is applied to each branch of the tree, and likelihood calculation of multiple alignments then becomes possible.

Since the TKF91 model is reversible, the root placement is immaterial. In the case of two sequences, one of the sequences may therefore be considered the ancestral sequence, and the other its descendant. In that framework, the number of possible alignments is finite. When the TKF91 model is applied to a tree, however, an infinite number of sequences may appear at internal nodes, for any given set of sequences which appear at the leaf nodes. Thus, such an extension is not straightforward, and several methods have been proposed. For star-shaped trees, Steel and Hein [SH01] have constructed an algorithm with  $O(L^{2k})$  running time for  $k$  sequences with geometric mean length  $L$ , and it has subsequently been extended, with similar time complexity, to binary trees [Hei01]. Thenceforth, the time complexity has been reduced to  $O(4^k L^k)$  for star-shaped trees [Mik02], and to  $O(5^k L^k)$  for arbitrary binary trees [HJP02] (see appendix). In the latter paper, the TKF91 model has been reformulated as a multiple-Hidden Markov Model, in which likelihood calculations can be performed using forward and backward algorithms known in the HMM literature [DEKM98]. Similar ideas have appeared for pairwise and triplewise statistical alignments as well [HB01]. In the Markov models there are  $O(\sqrt{5}^k)$  states, and therefore the running time of these algorithms is indeed  $O(5^k L^k)$ , in accordance with the general theory of HMMs [DEKM98]. The memory usage is of the order  $O(\sqrt{5}^k L^k)$  if the entire dynamic programming table is retained, and  $O(\sqrt{5}^k L^{k-1})$  if not [Hir77].

Because of the large number of states, the aforementioned algorithms are quite slow, even for moderate number of sequences. Hein *et al.* [HWK<sup>+</sup>00], however, have developed an algorithm for statistical alignment of two sequences which needs neither different states of a HMM nor partition of probabilities into different types of alignments [TKF91]. This is in contrast to the original formulation of the TKF91 model, which uses 3 states. The algorithm of [HWK<sup>+</sup>00] is, both in time and space complexity, as simple as the traditional distance- or similarity-based dynamic programming algorithms [NW70, SK83]. Henceforward, we refer to it as the *one-state recursion*.

In this paper, we present a one-state recursion for statistical alignment on arbitrary phylogenetic trees. Essentially, it combines the idea in [HWK<sup>+</sup>00] generalised to trees and Felsenstein's reverse traversal algorithm [Fel81]. The first idea allows us to reduce the number of states to one, and the second to compute in linear time the exponential number of terms that occur in the transition factors. The resulting algorithm has time complexity  $O(2^k L^k)$  and space complexity  $O(L^k)$  if the entire dynamic programming table is stored in

memory,  $O(L^{k-1})$  if not. This represents a great saving in both space and time compared to the hidden Markov recursion, and makes statistical alignment on a phylogenetic tree of modest size a definite practical possibility. Indeed, we have implemented both methods for 3 and 4 sequences, and achieved a considerable acceleration. One of us (IM) was able to perform a likelihood ratio test [Fel81, WLG01] for more than 250 triplets of yeast protein sequences in a day, using triplewise statistical alignments. Moreover, this one-state recursion can be coupled with corner-cutting techniques [HWK<sup>+</sup>00], which provides further reduction both in time and space complexity.

The organisation of this paper is as follows. We briefly describe the TKF model in §2 and discuss in the following section the one-state recursion in the case of two sequences. The two-sequence example is simple, but well illustrates our general approach. We discuss in §4 our algorithm for the one-state recursion and consider a specific application in §5. The main ideas underlying our proof are sketched in §6. The general one-state recursion is described in §7 and the computation of the transition factors which occur in the recursion is discussed in §8. In §9 we draw a connection with a hidden Markov model and describe how optimal alignments can be obtained. We conclude with some remarks and discussion in §10. In appendix, we show that the number of Markov states in the HJP recursion is  $O(\sqrt{5}^k)$  and give a proof of Lemma 1 from §7.

## 2 The TKF model

The TKF91 model is a continuous time reversible Markov model for the evolution of nucleotide (or amino acid) sequences. It models three of the main processes in sequence evolution, namely *substitutions*, *insertions* and *deletions* of nucleotides, approximating these as single-nucleotide processes. A nucleotide sequence is represented by an alternating string of *nucleotides* and *links*, connecting the nucleotides, and this string both starts and terminates with a link. We adopt the view that insertions originate from links, and add a nucleotide-link pair to the *right* of the original link; deletions originate from nucleotides and have the effect of removing the nucleotide and its right link. (This view is slightly different but equivalent to the original description, see [TKF91].) In this way, nucleotide subsequences evolve independently of each other, and the evolution of a nucleotide sequence is the sum of the evolutions of individual nucleotide-link pairs. The leftmost link of the sequence has no corresponding nucleotide to its left, hence it is never deleted, and for this reason it is called the *immortal link*.

Since subsequences evolve independently, it suffices to describe the evolution of a single nucleotide-link pair. In a given time-span  $\tau$ , this evolves into a sequence of nucleotides of finite length. The first nucleotide of this sequence may be homologous to the original one, while subsequent ones will be non-homologous. Table 1 summarises the corresponding probabilities. The parameters  $\lambda$  and  $\mu$  are the birth rate per link and the death rate per nucleotide, respectively, and in order to have a finite equilibrium sequence length, we require

Fate:	Probability:	Label:
$C \rightarrow C\#^{n-1}$	$e^{-\mu\tau}(1 - \lambda\beta(\tau))(\lambda\beta(\tau))^{n-1}$	$H_\tau B_\tau^{n-1}$
$C \rightarrow \#^n$	$(1 - e^{-\mu\tau} - \mu\beta(\tau))(1 - \lambda\beta(\tau))(\lambda\beta(\tau))^{n-1}$	$N_\tau B_\tau^{n-1}$
$C \rightarrow -$	$\mu\beta(\tau)$	$E_\tau$
$\star \rightarrow \star\#^n$	$(1 - \lambda\beta(\tau))(\lambda\beta(\tau))^n$	$I_\tau B_\tau^n$

Table 1: Probabilities of the TKF91 model (see text).

$\lambda < \mu$ . For brevity we write

$$\beta(\tau) = \frac{1 - e^{(\lambda-\mu)\tau}}{\mu - \lambda e^{(\lambda-\mu)\tau}}.$$

On the right-hand side of the arrow in the column labeled “Fate”,  $C$  denotes a nucleotide homologous to the original nucleotide, whereas  $\#$ ’s denote non-homologous nucleotides. The immortal link is denoted by  $\star$  and other links are suppressed. All final arrangements can be thought of as being built from five basic “processes” which we call Birth, Extinction, Homologous, New (or Non-homologous) and Initial (or Immortal). These processes are labeled by their initials, and each corresponds to a specific probability factor as follows:

$$\begin{aligned} B_\tau &= \lambda\beta(\tau) & E_\tau &= \mu\beta(\tau) & I_\tau &= 1 - \lambda\beta(\tau) \\ H_\tau &= e^{-\mu\tau}(1 - \lambda\beta(\tau)) & N_\tau &= (1 - e^{-\mu\tau} - \mu\beta(\tau))(1 - \lambda\beta(\tau)) \end{aligned} \quad (1)$$

In a tree, time flows forward from the root to the leaves, and to each node of the tree we associate a time parameter  $\tau$  which is set equal to the length of the incoming branch. For the root,  $\tau = \infty$ , and the resulting equilibrium length distribution of the immortal link sequence is geometric with parameter  $B_\infty = \lambda/\mu$ ; other links will have left no descendants since  $H_\infty = N_\infty = 0$ .

In the original TKF91 model, a simple substitution process known as the “Felsenstein81 model” [Fel81] has been used. It is straightforward to replace this by more general models for substitutions of nucleotides or amino acids [HWK<sup>+</sup>00]. In the present paper, when a new non-homologous character appears at a node (as the result of a  $B$  or  $N$  process), it is always drawn from the equilibrium distribution; if a character at a node is homologous to the character at its immediate ancestral node, then the probability of this event is given by the chosen substitution model.

Note that we do not sum over possible *alignments*. Although for two sequences there almost<sup>1</sup> is a correspondence between alignments and evolutionary history (in the TKF model), this correspondence breaks down for trees with larger number of leaves, since information on what occurs at the internal nodes is lost in an alignment. The homology structure of nucleotides is the only part of the evolutionary history relevant to the alignment.

---

<sup>1</sup>As was pointed out in [TKF91], gaps in two sequences that immediately follow each other may be interchanged without changing the interpretation as alignments, however the evolutionary interpretation in the TKF model is different.

### 3 One-state recursion: A graphical proof for two sequences

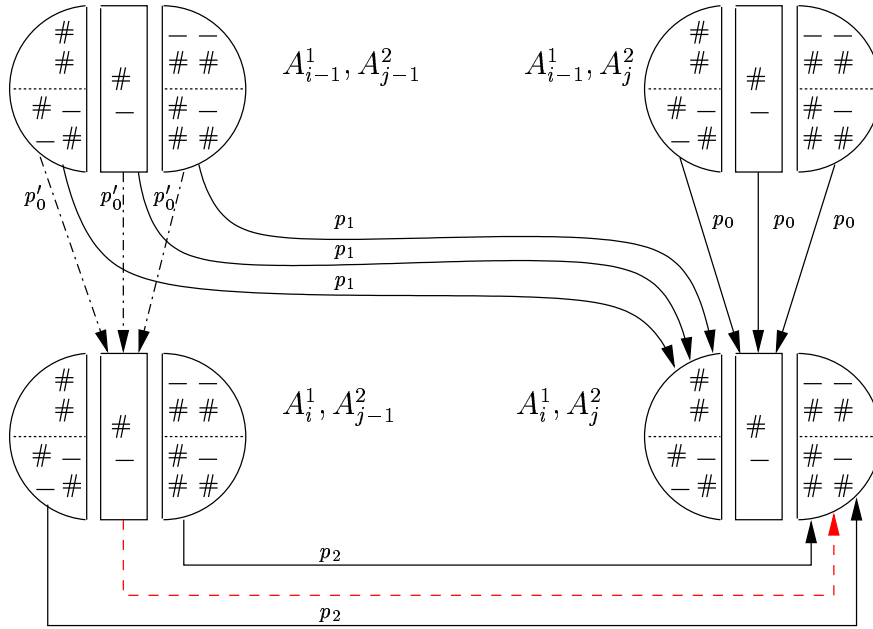


Figure 1: A graphical representation of the three-state recursion for two-sequence alignment. The  $p_k$  are transition probabilities associated to  $S^k$  (see text):  $p_0 = B_\infty E_\tau$ ,  $p_1 = B_\infty H_\tau + B_\infty N_\tau$  and  $p_2 = B_\tau$ .

In [HWK<sup>+</sup>00] a proof has been given for the one-state recursion in the case of two-sequence alignment. Here we provide a different proof for this particular case. Our proof for the general case follows a similar line of ideas.

Suppose  $A^1$  is the ancestral sequence observed at the root and  $A^2$  is the descendant sequence which results from the TKF91 process after time  $\tau$ . Let  $A_i^k$  (resp.  $a_i^k$ ) denote the  $i$ -long prefix (resp. the  $i$ th character) of sequence  $A^k$ . Let  $P(i, j)$  be the joint probability of observing  $A_i^1$  and  $A_j^2$ . Given that  $P(i-1, j)$ ,  $P(i, j-1)$  and  $P(i-1, j-1)$  are known, we want to compute  $P(i, j)$ .

The possible evolutionary histories (or alignments) are customarily [TKF91] classified into three groups,  $S^0$ ,  $S^1$  and  $S^2$ , according to whether the rightmost root link has 0, 1 or at least 2 descendants at the leaf. Another way to define these sets is by the last columns in the alignment; symbolically,  $S^0 = \{ \# \_ \}$ ,  $S^1 = \{ \# \_ \_ \}$  and  $S^2 = \{ \# \_ \_ \_ \}$ , where the upper symbols represent nucleotides at the root. These sets are associated to the evolutionary processes as follows. After the  $B_\infty E_\tau$  process, the rightmost root link has 0 descendants, so the history is in  $\{ \# \_ \}$ . After  $B_\infty H_\tau$  or  $B_\infty N_\tau$  it is in  $\{ \# \_ \_ \}$ , while after a birth at the leaf node ( $B_\tau$ ), the alignment ends in either of  $\{ \# \_ \_ \_ \}$ . These processes may happen in any order, except that a transition from  $\{ \# \_ \}$  to  $\{ \# \_ \_ \_ \}$  is not allowed: new births at the leaf ( $B_\tau$  processes) are allowed except after an extinction ( $E_\tau$ ) process, see table 1. A graphical depiction of the three-state recursion that implements this restriction is provided



## 4 Algorithm

Using the TKF model, it is easy to compute the likelihood of a specific evolutionary history, given as a sequence of events. In effect, this integrates out all possible evolutionary histories on the edges of the phylogenetic tree, while fixing the nucleotide sequences and their mutual alignment at internal nodes. The likelihood of observing the given sequences at the leaves is obtained by summing this likelihood over all possible sequences of events. Since we usually do not have information about the sequences at the internal nodes, such a computation should sum out the internal sequences as well. This is what our algorithm computes.

In order to write down the main theorem, we need some notation.  $T$  is the set of nodes of the phylogenetic tree relating the sequences, and  $r$  is its root. If  $n \in T$  is a node, then  $n_l$  and  $n_r$  denote the left and right immediate descendants of  $n$ . For a node  $n$ ,  $\tau(n)$  denotes the length of the incoming branch, and we abbreviate  $B_n := B_{\tau(n)}$  etc.. For the root,  $\tau(r) = \infty$  by definition. Finally,  $p_n(\alpha \rightarrow \gamma)$  is the probability that nucleotide  $\alpha$  evolves into  $\gamma$  in time  $\tau(n)$ .

**Theorem 1** *Suppose  $A^1, \dots, A^k$  are sequences related by a phylogenetic tree  $T$  with  $k$  leaves and root  $r$ . Let  $\mathbf{K} = (K_1, \dots, K_k)$  and let  $P(\mathbf{K})$  be the probability of producing, under the TKF91 model, the prefixes of  $A^1, \dots, A^k$  up to position  $K_1, \dots, K_k$ , respectively, where the numbering starts at 0 with the immortal link. The following algorithm computes  $P(\mathbf{K})$ :*

$$P(\mathbf{0}) = \frac{\prod_{n \in T} (1 - B_n)}{G_{\mathbf{0}}^{\mathbf{0}}(r, -)},$$

$$P(\mathbf{K}) = \frac{1}{G_{\mathbf{K}}^{\mathbf{0}}(r, -)} \sum_{\mathbf{v} \in \{0,1\}^k \setminus \mathbf{0}} (-G_{\mathbf{K}}^{\mathbf{v}}(r, -)) P(\mathbf{K} - \mathbf{v}).$$

If  $n$  is an internal node, then

$$G_{\mathbf{K}}^{\mathbf{v}}(n, \alpha) = \left[ E_{n_l} G_{\mathbf{K}}^{\mathbf{v}}(n_l, -) + \sum_{\gamma} \left( H_{n_l} p_{n_l}(\alpha \rightarrow \gamma) + N_{n_l} \pi(\gamma) \right) G_{\mathbf{K}}^{\mathbf{v}}(n_l, \gamma) \right]$$

$$\times \left[ E_{n_r} G_{\mathbf{K}}^{\mathbf{v}}(n_r, -) + \sum_{\gamma} \left( H_{n_r} p_{n_r}(\alpha \rightarrow \gamma) + N_{n_r} \pi(\gamma) \right) G_{\mathbf{K}}^{\mathbf{v}}(n_r, \gamma) \right],$$

$$G_{\mathbf{K}}^{\mathbf{v}}(n, -) = G_{\mathbf{K}}^{\mathbf{v}}(n_l, -) G_{\mathbf{K}}^{\mathbf{v}}(n_r, -) - B_n \sum_{\gamma} \pi(\gamma) G_{\mathbf{K}}^{\mathbf{v}}(n, \gamma).$$

If  $n$  is a leaf node, then

$$G_{\mathbf{K}}^{\mathbf{v}}(n, \alpha) = \begin{cases} 1, & \text{if } v_n = 1 \text{ and } a_{K_{n-1}}^n = \alpha, \\ 0, & \text{otherwise.} \end{cases}$$

$$G_{\mathbf{K}}^{\mathbf{v}}(n, -) = \begin{cases} 1, & \text{if } v_n = 0, \\ -B_n \pi(a_{K_{n-1}}^n), & \text{if } v_n = 1. \end{cases}$$

By  $K_n$  we mean the component of the vector  $\mathbf{K}$  corresponding to leaf  $n$ , and similarly for  $v_n$ . By  $a_{K_n-1}^n$  we denote the character at position  $K_n - 1$  in the sequence at leaf  $n$ .

The numbers  $G_{\mathbf{K}}^{\mathbf{v}}(r, -)$  are the *transition factors* and play the same role as the transition factors in Figure 2. They are sums and differences of probabilities, designed to exactly cancel all contributions of illegal transitions due to the merging of states into a single probability, analogous to the two-sequence case described in §3.

Note that  $G_{\mathbf{K}}^{\mathbf{0}}$  does not depend on  $\mathbf{K}$ , so a small acceleration can be achieved by computing the pre-factor  $1/G_{\mathbf{K}}^{\mathbf{0}}$  only once. A more substantial acceleration is achieved by noting that  $G_{\mathbf{K}}^{\mathbf{v}}(n, \alpha)$  and  $G_{\mathbf{K}}^{\mathbf{v}}(n, -)$  depend only on the components of  $\mathbf{v}$  which correspond to the leaves contained in the subtree of  $n$ . Therefore, if the tree with root  $n$  has  $s$  leaves, it is enough to compute  $G_{\mathbf{K}}^{\mathbf{v}}(n, \alpha)$  for  $2^s$  values of  $\mathbf{v}$ . Once these are computed, any entry for the ancestral node can be computed in constant time. Therefore the overall running time of the algorithm is  $2^k L^k$ , where  $k$  is the number of sequences and  $L$  is the geometric mean of the sequence lengths.

## 5 Example

As a concrete example, we have used our algorithm to find the most likely evolutionary tree which relates the following four globins: human  $\alpha 1$  and  $\beta$  hemoglobins, human myoglobin, and leghemoglobin from the bean *Canavalia lineata*. There are three topologically different unrooted trees with four leaves, and we have estimated maximum likelihood parameters for each of them. The most likely tree is, as expected, the one that groups human  $\alpha 1$  and  $\beta$  hemoglobin together; the likelihoods of the other two trees are about a factor 400 smaller (see Figure 3). After obtaining the maximum likelihood tree, the maximum likelihood alignment can be obtained using the Viterbi algorithm. This alignment is shown in Table 2. For a more detailed discussion on the Viterbi algorithm, see §9.

```
Hba1: MV--LSPADKTNVKAAWGKVGAHAGEYGAELERMFLSFPTTKTYFPHF--DLS-H-----GSAQVKGHGKKVAD-AL-TNA-
Hbb:  MV-HLTPEEKSAVTALWGKV--NVDEVGGEALGRLLVVYPWTQRFFESF-GDLSTPDAVM-GNPKVKAHGKKVLG-AF-SDG-
Myo:  MG--LSDGEWQLVLNVWGKVEADIPGHGQEVLIIRLFKGGHPETLEKFDKFK-HLKSEDE-MKASEDLKKHGATVLT-AL-GGI-
Legh: MGA-FSEKQESLVKSSWEAFKQNVPHHSAVFYTLILEKAPAAQNMFS-F---LSNGVD-P-NNPKLKAHAEKVFKMTVDSAVQ

      VAHVDDMPNALSALSDLHAHKLRVDPVNFK-LLSHCLLVTLAAHLPAEFTPAVHASLDKFLASVSTVL-TS-K---YR-
      LAHLDNLKGTFAATLSELHCDKLHVDPENFR-LLGNVLVCLAHHFGKEFTPPVQAAAYQKVVAGVANAL-AH-K---YH-
      LKKKGHHEAEIKPLAQSHATKHKI-PVKYLEFISECIIQVLQSKHPGDFGADAQGAMNKALELFRKDMASNYKELGFQG
      LRAKGEVVLADPTLGSVHVQKGVLDP-HFL-VVKEALLKTFKEAVGDKWDELGNAWEVAYDELAALAI-KK-A-MGSA-
```

Table 2: The maximum likelihood alignment for the first pedigree in Figure 3. The log-likelihood of this alignment is  $-1593.223$ .



We have optimised the likelihood as a function of the five branch length parameters and  $\lambda$ ; the deletion rate parameter  $\mu$  has been adjusted to make the expected sequence length  $\frac{\lambda}{\mu-\lambda}$  equal to the average sequence length. By using corner cutting methods we could confine the calculation to the parts of the dynamic programming table which contribute non-infinitesimally to the final probability, resulting in a speedup of about a factor of 500. At the extremal point, an estimate of the second derivative yielded a covariance matrix; see Figure 3 for the standard deviations.

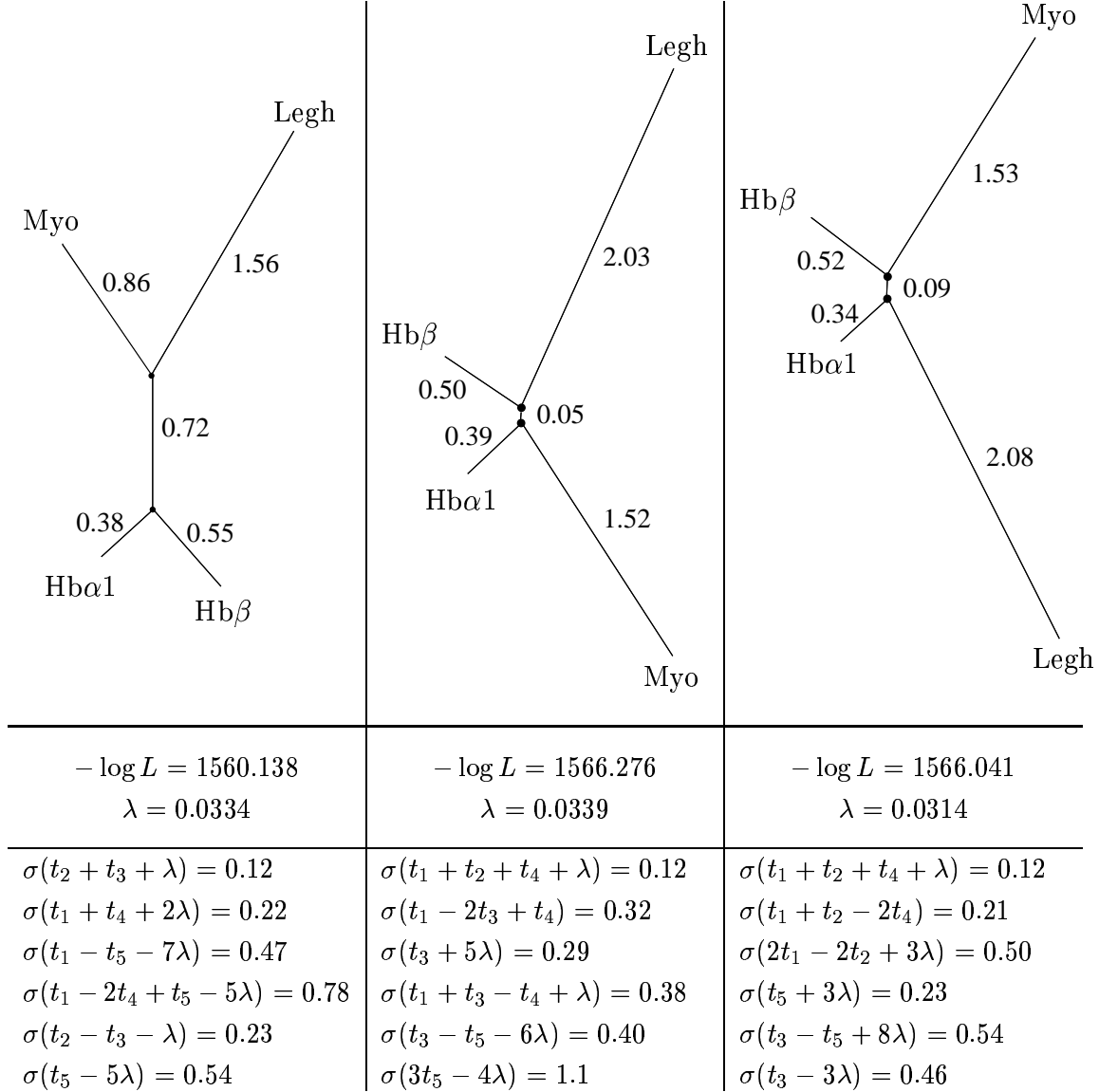


Figure 3: Maximum likelihood trees relating human  $\alpha 1$  and  $\beta$  hemoglobins, myoglobin, and bean leghemoglobin, for all three topologically distinct trees; total likelihood values and insertion rate; and estimated standard deviations. The numbers refer to the branch length in units of expected number of substitutions per amino acid.

## 6 Idea of proof

Our aim is to give a “one-state” dynamic programming algorithm for calculating the joint likelihood of observing a set of sequences at the leaves of an evolutionary tree. This likelihood is the sum of the probabilities of all evolutionary histories which produce the observed sequences at the leaves.

We introduce the concept of an “event” which corresponds to a nucleotide birth at a node of the phylogenetic tree and its subsequent fate down its descendant subtree. More precisely, in the process of going from one node to a descendant node, one of the following three things may happen: The nucleotide may survive as a homologous nucleotide, it may die leaving at least one surviving new nucleotide, or it may go extinct altogether. These three possibilities are labeled H, N and E respectively, as discussed in §2. In the first two cases, other non-homologous descendants may have been born, but these possibilities are dealt with in subsequent events. Hence, in this model, evolutionary histories of nucleotide sequences correspond to sequences of events. However this correspondence is not one-to-one for two reasons. One reason is that we need to honour the TKF model, which does not allow extinct nucleotides to spawn new ones. (In Table 1, this is reflected by the third entry,  $C \rightarrow -$ , which is the only one that may not be followed by a string of #’s.) The other reason is that different sequences of events may represent the same evolutionary history, since the relative ordering of events on disjoint subtrees has no meaning. In a Markov chain approach, these two requirements are met by defining states and by ruling out certain events (*i.e.* transitions, in the Mealy machine view) in certain states. Essentially that is how we also approach the problem. We define a “state” and use it to rule out events which violate the TKF model. Our choice of “state” still allows over-counting of histories. To overcome this, we require the events to be *ordered* in a specific way depending on the state.

The ordering we use has the property that in properly ordered sequences of events, an event which violates the TKF model can be recognised by comparing it to its immediate predecessor. This means that we do not need the state to decide which events are allowed, but the state is still used to define the proper ordering. We can now use a recursion which, as a first approximation, allows all events. This includes all *legal* sequences of events, as well as some sequences of events which end in an illegal pair, either because it is not properly ordered, or because it violates the TKF model. As a second approximation, we subtract all sequences of events which end in particular illegal pairs. In the same way as before, this includes some illegal triplets which should not have been subtracted, and these are added in again, etc., in a procedure reminiscent of inclusion-exclusion. We have dubbed the pairs, triplets etc. involved in this procedure the *disallowed sequences*.

It turns out that there exist only finitely many disallowed sequences, so that the inclusion-exclusion procedure stops after finitely many steps (bounded by the number of nodes in the tree). Moreover, disallowed sequences emit at most one nucleotide per leaf, and that leads to a recursion which refers back only to the  $2^k - 1$  nearest predecessors in the  $k$ -dimensional

dynamic programming table, where  $k$  is the number of leaves.

The ordering, and hence the disallowed sequences we have to sum over, still depends on the state. However, the state turns out to influence only the *ordering* of the elements of these disallowed sequences, not the set of events themselves. More precisely, for each state, there is a one-to-one correspondence between such sets (that we call *illegal sets*) and disallowed sequences. The terms involved in the summation do not depend on the ordering, and this enables us to sum away the dependence on the state. A small technical problem is that some events do not emit any character (known in the HMM literature as “silent states”). These events lead to self-references in the recursion, turning it into a (trivial, one-dimensional) linear equation. We can sum over infinitely many possible silent events using the standard “wing folding” technique [Edd01], which amounts to solving the linear equation.

It remains to compute the transition factors, which involve sums over illegal sets of events that lead to specific emissions. The number of illegal sets is still exponential in the number of nodes in the phylogenetic tree. However, using a correspondence between illegal sets and labelings of the tree, it is possible to compute this sum in linear time using a reverse-traversal algorithm. This recursion also deals with the nucleotide assignments in the internal nodes of the tree, similar to Felsenstein’s postorder tree traversal algorithm [Fel81].

## 7 The general one-state recursion

In this section we introduce the general one-state recursion. The main result is presented in §7.2; to maintain the flow of discussion, we defer the proof until Appendix B. In §8, to arrive at the algorithm discussed in §4, we combine the result of the present section with a fast algorithm for computing the transition factors.

### 7.1 Definitions

Let  $T$  denote the set of nodes of a phylogenetic binary tree. Our result can be extended to more general trees, but that is not pursued here. Note that we suppose the tree edges to be given and fixed, and that we only include the nodes in our representation  $T$ . A *subtree*  $t$  of  $T$  is a set of nodes with the properties that (i) for any node  $n \in t$ , its descendants are in  $t$ , and (ii) there is a unique node, called the *root* of  $t$  and denoted by  $r(t)$ , without a parent. A *labeling* of a tree  $t$  is a function  $l : t \rightarrow \{B^\alpha, -, H^\alpha, N^\alpha \mid \alpha \in \mathcal{A}\}$ , where  $\mathcal{A}$  is the set alphabet, which is  $\{A, C, G, T\}$  in the case of DNA sequences.  $B$ ,  $H$  and  $N$  stand for birth, homologous and new, respectively, as discussed in §2, and  $\alpha$  is the character associated to the process. The symbol “ $-$ ” stands for “no character” and means that an  $E_\tau$  (extinction) process has occurred somewhere higher up in the tree.

An *event*  $e$  is a pair  $(e^t, e^l)$ , where  $e^t$  is a subtree of  $T$  and  $e^l$  a labeling of  $e^t$ , with the additional properties that (i)  $e^l(n) = B^\alpha$  (any  $\alpha$ ) if and only if  $n = r(e^t)$ , and (ii) if

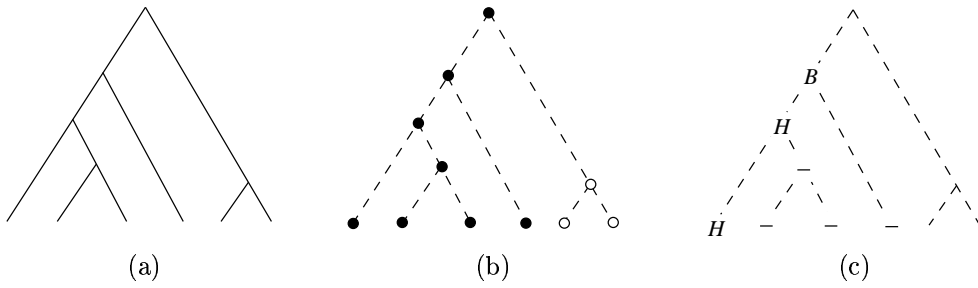


Figure 4: (a) A 6-leaved phylogenetic binary tree. (b) A possible state. Filled circles ● (resp. open circles ○) represent the nodes which are contained (resp. not contained) in the state. (c) A possible event. Any node without a label is not contained in the subtree corresponding to the event.

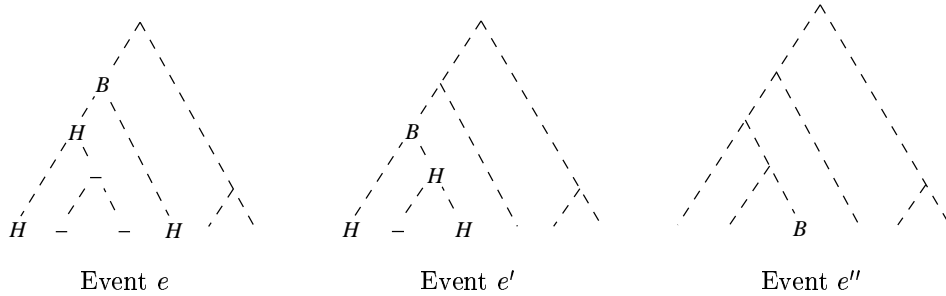


Figure 5: Three possible events.

$e^l(n) = -$  for any  $n \in e^t$ , then  $e^l(n') = -$  for all descendants  $n'$  of  $n$ . In this way, an event represents the birth of a nucleotide at  $r(e^t)$  and its subsequent fate down the phylogenetic tree. An event carries a definite probability  $p(e)$  which can be calculated from the TKF91 model, and we will concern ourselves with that in §8. Note that an event gives rise to at most one new nucleotide at each leaf sequence, and that such nucleotides may or may not be homologous to each other.

We consider *sequences of events*, denoted by  $E = (e_1, e_2, \dots, e_N)$ . These give rise to sequences of *states*  $(S_1, S_2, \dots, S_N)$ , which are defined recursively as follows:  $S_1 := T$  and  $S_{i+1} := (S_i \setminus e_i^t) \cup \{n | n \in e_i^t \text{ and } e_i^l(n) \neq -\}$ . The state  $S_i$  determines at which nodes new nucleotides may be born at event  $i$ . Figure 4 illustrates these definitions.

We introduce some notation for convenience:  $t_i := e_i^t$ ,  $r_i := r(e_i^t)$ ,  $l_i := e_i^l$ . These stand for the *tree*, the *root*, and the *labeling* of event  $e_i$ . We define the *context*  $c_i$  by  $c_i := F$  if  $r_i \in S_i$ , and  $c_i := I$  if  $r_i \notin S_i$ , meaning Fertile if the birth of event  $e_i$  occurred at a live node, Infertile if not. An event  $e_i$  is called *legal* if  $c_i = F$ . A sequence of events is called *legal* if all events are legal. A sequence of legal events results in a *legal state*, which equivalently is a state that can be reached by a single event at the root. Examples of legal and illegal states are provided in Figure 6.

In a sequence  $E$  of events, the ordering of the events is important insofar as it determines the ordering of events pertaining to each individual leaf sequence. Permutations do not

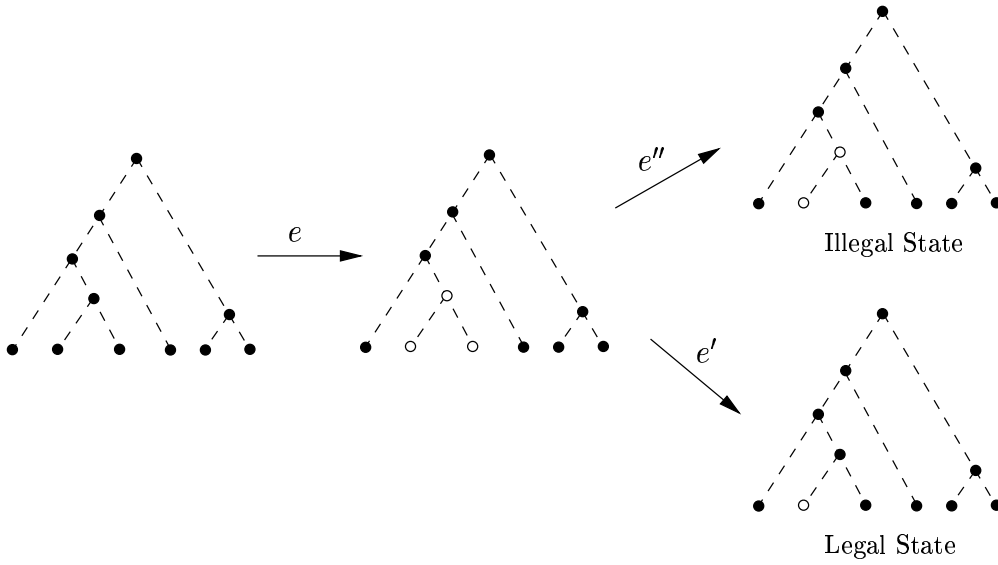


Figure 6: Sequences of states corresponding to the sequences  $(e, e')$  and  $(e, e'')$  of events, where  $e, e', e''$  are as in Figure 5. The sequence  $(e, e')$  is a legal sequence, resulting in a legal state, whereas  $(e, e'')$  is an illegal sequence, giving rise to an illegal state. Note that  $(e, e', e'')$  is a legal sequence.

change the evolutionary history represented by a sequence as long as, for all leaves  $n$ , the ordering in the subsequence of those events which intersect the path from root to  $n$  is undisturbed. More precisely, two sequences  $E = (e_1, \dots, e_N), E' = (e'_1, \dots, e'_N)$  represent the same history iff there exists a permutation  $\sigma$  of  $\{1, \dots, N\}$  such that

$$e_{\sigma(i)} = e'_i \quad \text{and} \quad i < j, \quad t_i \cap t_j \neq \emptyset \Rightarrow \sigma(i) < \sigma(j).$$

We call these  $\sigma$  *allowed permutations* (for  $E$ ). This defines an equivalence relation on the set of sequences of events. We call the equivalence class of any sequence of events the (*evolutionary*) *history* defined by that sequence.

Although in general the state depends on the ordering of events, the context  $c_j$  depends only on the history, since it depends on whether  $r_j \in S_j$ , which in turn depends on  $l_i(r_j)$ , for the unique  $i < j$  with  $r_j \in t_i$  and  $r_j \notin t_k$  for all  $i < k < j$ . The pair  $(j, i)$  so defined is equivariant under allowed permutations, *i.e.* the indices of the pair transform according to the permutation. This allows us to define a *legal history* to be a history represented by a legal sequence of events.

## 7.2 Summing away the states

We denote by  $E(\mathbf{K})$  a set of legal sequences of events which emit the prefixes of the given sequences  $A^1, \dots, A^k$  up to position  $\mathbf{K} = (K_1, \dots, K_k)$ , and which contains precisely one

representative of each history. The quantity of interest is the likelihood

$$P(\mathbf{K}) := \sum_{(e_1, \dots, e_N) \in E(\mathbf{K})} p(e_1) \cdots p(e_N).$$

To state the main result of this section, we need to introduce one more definition.

**Definition 1 (Illegal Set)** *A set of events  $\{e_1, \dots, e_N\}$  is called an illegal set if  $t_i \supseteq t_j \Rightarrow l_i(r_j) = -$ .*

**Lemma 1** *The likelihood  $P(\mathbf{K})$  satisfies the following equation:*

$$\begin{aligned} P(\mathbf{K}) &= \sum_e P(\mathbf{K} - \mathbf{v}_e) p(e) \\ &\quad - \sum_{\{e_1, e_2\} \text{ illegal}} P(\mathbf{K} - \mathbf{v}_{e_1} - \mathbf{v}_{e_2}) p(e_1) p(e_2) \\ &\quad + \sum_{\{e_1, e_2, e_3\} \text{ illegal}} P(\mathbf{K} - \mathbf{v}_{e_1} - \mathbf{v}_{e_2} - \mathbf{v}_{e_3}) p(e_1) p(e_2) p(e_3) \\ &\quad - \dots \end{aligned} \tag{2}$$

Here  $\mathbf{v}_e$  is the emission vector corresponding to the event  $e$ , that is, a  $k$ -dimensional vector whose component corresponding to each leaf is 0 if the leaf is labeled  $-$ , or 1 otherwise.

**Proof:** See Appendix B.

### 7.3 Eliminating the silent events

Note that, for a nontrivial tree  $T$ , the set  $E(\mathbf{K})$  of legal sequences of event emitting a certain prefix is infinite, even for  $\mathbf{K} = \mathbf{0}$ . This is due to the so-called silent events, *i.e.* events without emissions at the leaves. This results in the appearance of  $P(\mathbf{K})$  on the right-hand side of (2). However, the contribution of these silent events can easily be calculated using a trick analogous to the “wing folding” technique in the HMM literature [Edd01]. A similar trick can be found in [SH01].

Note that, by definition,  $\{e\}$  is an illegal set for all events  $e$ , and therefore (2) can be rewritten as follows:

$$P(\mathbf{K}) = \sum_{M \text{ illegal set,}} P(\mathbf{K} - \sum_{e_i \in M} \mathbf{v}_{e_i}) (-1)^{|M|+1} \prod_{e_i \in M} p(e_i). \tag{3}$$

We can decompose the sum as

$$\begin{aligned} P(\mathbf{K}) &= \sum_{\substack{M \text{ illegal set,} \\ \text{not all } \mathbf{v}_{e_i} = \mathbf{0}}} P(\mathbf{K} - \sum_{e_i \in M} \mathbf{v}_{e_i}) (-1)^{|M|+1} \prod_{e_i \in M} p(e_i) \\ &\quad + P(\mathbf{K}) \sum_{\substack{M \text{ illegal set,} \\ \text{all } \mathbf{v}_{e_i} = \mathbf{0}}} (-1)^{|M|+1} \prod_{e_i \in M} p(e_i). \end{aligned}$$

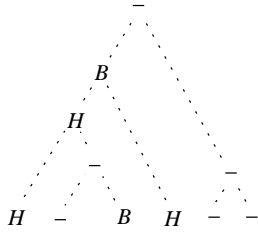


Figure 7: The tree labeling corresponding to the illegal set  $\{e, e''\}$ , where the events  $e$  and  $e''$  are as in Figure 5.

Solving this equation for  $P(\mathbf{K})$  gives the correct answer and includes histories with  $0, 1, 2, \dots$  silent events. A heuristic way to see this is to observe that the equation  $X = a + sX$ , where  $s$  and  $a$  are the silent and non-silent contributions, respectively, is solved by  $X = a/(1 - s) = a(1 + s + s^2 + \dots)$ . This expansion clearly shows the separate contributions of chains of silent events of length  $0, 1, 2, \dots$ . The reasoning in the proof of Lemma 1 makes sure that only sequences of legal events are included, and that applies to the silent states as well.

## 8 Pruning the tree

The recursion (2) enables us to compute  $P(\mathbf{K})$ , but it requires summing over a number of terms exponential in the number of sequences. A similar problem occurs when calculating the likelihood of nucleotide emissions on a tree under a simple substitution model, in which case one needs to sum over an exponential number of nucleotide assignments to internal nodes. Dynamic programming helps out once again, in the form of Felsenstein’s linear time post-order tree traversal algorithm. We proceed in a similar way, by finding a correspondence between certain labelings of  $T$  and illegal sets of events. The summation over all such labelings, including the  $(-1)^{k+1}$  sign arising from the inclusion-exclusion argument, can be performed in linear time by a dynamic programming algorithm similar to Felsenstein’s.

The probability of an event is the product of conditional probabilities at the nodes of the phylogenetic tree, and the conditional probability at a node only depends on the labeling of the node and its ancestral node. The label determines both the process and the associated nucleotide, except for “-”, which implies an Extinction event only if its ancestral node is not labeled “-”. Symbolically,

$$\begin{aligned}
 p(N^\alpha | X^\gamma) &= N_n \pi(\alpha), & p(B^\alpha | -) &= B_n \pi(\alpha), \\
 p(H^\alpha | X^\gamma) &= H_n p_n(\gamma \rightarrow \alpha), & p(N^\alpha | -) &= 0, \\
 p(- | X^\gamma) &= E_n, & p(H^\alpha | -) &= 0, \\
 & & p(- | -) &= 1,
 \end{aligned}$$

where  $\alpha, \gamma \in \mathcal{A}$  (c.f. §7.1). The symbol  $X$  may be anything except “-”. The probability factors  $B_n, N_n, H_n$  and  $E_n$  are subscripted with a node, and are defined (c.f. (1)) in terms

of the length of the incoming branch ( $\infty$  if  $n = r(T)$ ). Note that  $p(B^\alpha|X^\alpha)$  will never occur.

This allows us to calculate probabilities of events, but we would like to sum over illegal sets directly. It follows from Definition 1 that for each node  $n \in T$ , at most one event in the illegal set has  $n$  labeled with a symbol other than “-”. Furthermore, since the roots of events are uniquely identified as the only nodes labeled  $B^\alpha$ , there is a one-to-one correspondence between illegal sets of events and labelings of a single tree with  $B^\alpha, H^\alpha, N^\alpha, -$  which obey the following rule: the root and the immediate descendants of “-” are labeled with  $B^\alpha$  or “-”. See Figure 7 for an example. Indeed, such labelings can be decomposed into illegal sets of events by recursively removing subtrees containing exactly one  $B^\alpha$  label at their roots, and resetting the labelings to “-”. Note that the all-“-” tree has a proper labeling according to these rules, and corresponds to the empty illegal set.

The factor associated to an illegal set  $S$  is the product of the separate event probabilities, with a sign  $(-1)^{|S|+1}$ . The exponent  $|S| + 1$  is just one more than the number of  $B^\alpha$ 's occurring in the associated tree labeling. Hence, analogous to computing the probability of an event, the factor associated to a tree labeling may be calculated (except for a single overall minus sign) as the product of several “conditional factors” as follows:

$$\begin{aligned} f(B^\alpha|X^\gamma) &= -E_n B_n \pi(\alpha), & f(B^\alpha|-) &= -B_n \pi(\alpha), \\ f(N^\alpha|X^\gamma) &= N_n \pi(\alpha), & f(N^\alpha|-) &= 0, \\ f(H^\alpha|X^\gamma) &= H_n p_n(\gamma \rightarrow \alpha), & f(H^\alpha|-) &= 0, \\ f(-|X^\gamma) &= E_n, & f(-|-) &= 1. \end{aligned}$$

If we group terms in (3) according to their total emission vector  $\mathbf{v} = \mathbf{v}_{e_1} + \dots + \mathbf{v}_{e_m}$ , the sum of contributing factors may be obtained by a Felsenstein-like recursion. Since the conditional factors above only depend on  $\alpha$  and on whether the label is “-”, the following recursion computes the required sum:

$$\begin{aligned} F(n, \alpha) &= \left[ E_{n_l} F_{\mathbf{K}}^{\mathbf{v}}(n_l, -) + \sum_{\gamma} F(n_l, \gamma) (H_{n_l} p_{n_l}(\alpha \rightarrow \gamma) + N_{n_l} \pi(\gamma) - E_{n_l} B_{n_l} \pi(\gamma)) \right] \\ &\quad \times \left[ E_{n_r} F(n_r, -) + \sum_{\gamma} F(n_r, \gamma) (H_{n_r} p_{n_r}(\alpha \rightarrow \gamma) + N_{n_r} \pi(\gamma) - E_{n_r} B_{n_r} \pi(\gamma)) \right], \\ F(n, -) &= \left[ F^{\mathbf{v}}(n_l, -) - \sum_{\gamma} B_{n_l} F^{\mathbf{v}}(n_l, \gamma) \pi(\gamma) \right] \times \left[ F^{\mathbf{v}}(n_r, -) - \sum_{\gamma} B_{n_r} F^{\mathbf{v}}(n_r, \gamma) \pi(\gamma) \right]. \end{aligned}$$

Here,  $F(n, \alpha)$  (resp.  $F(n, -)$ ) is the sum of all products of conditional factors on the subtree with root  $n$ , given that  $n$  is labeled  $X^\alpha$  (resp. “-”).

Note that we suppressed the dependence on  $\mathbf{v}$ . This dependence only surfaces when  $n$  is a leaf, in which case the recursion terminates with  $F(n, \alpha) = 1$  if  $\mathbf{v}_n = 1$  and  $\alpha$  is the character in sequence  $A^n$  at the current position, and  $F(n, -) = 1$  if  $\mathbf{v}_n = 0$ ;  $F = 0$  in all



other cases. The transition factor associated to  $\mathbf{v}$  is finally

$$\sum_{\substack{\{e_1, \dots, e_k\} \text{ illegal} \\ \mathbf{v}_{e_1} + \dots + \mathbf{v}_{e_k} = \mathbf{v}}} (-1)^{k+1} \prod_i p(e_i) = - \left( F(r, -) - \sum_{\alpha} F(r, \alpha) B_r \pi(\alpha) \right),$$

where  $r$  is the root of  $T$ . A final useful simplification occurs if we set  $G(n, \alpha) = F(n, \alpha)$  and  $G(n, -) = F(n, -) - \sum_{\gamma} B_n \pi(\gamma) F(n, \gamma)$ , in which case the recursion for  $G$  takes the form given in §4 and the transition factor is simply  $-G(r, -)$ .

To compute  $P(\mathbf{K})$  we have to eliminate the silent events by solving a linear equation, as described in §7.3. This amounts to dividing each transition factor by

$$1 - \sum_{\substack{\{e_1, \dots, e_k\} \text{ illegal} \\ \mathbf{v}_{e_1} + \dots + \mathbf{v}_{e_k} = \mathbf{0}}} (-1)^{k+1} \prod_i p(e_i), \quad (4)$$

one minus the transition factor associated to the null emission. Since the recursion for  $F$  also includes the “-”-labeled tree, which gets assigned a term  $+1$ , in fact  $G(r, -)$  for  $\mathbf{v} = \mathbf{0}$  is precisely equal to (4).

It remains to compute the initial value  $P(\mathbf{0})$ , the probability of emitting no nucleotides to any sequence and being left with only the immortal link. According to the TKF91 model, this is the product of the immortal link pre-factors  $I_n = 1 - B_n$ , where the product extends over all nodes in  $T$ , and silent events are included by dividing this by  $G(r, -)$ . This completes the proof of Theorem 1. ■

## 9 A Hidden Markov model

The algorithm of Theorem 1 efficiently computes the likelihood of aligning sequences in a tree and may be used to find maximum likelihood parameters, but it cannot be used to find the most likely alignment. To obtain the optimal alignment, we need to apply the Viterbi algorithm to the full hidden Markov model (HMM). We define such an HMM in this section. The one-state recursion is still useful for defining a narrow “region of interest” in which one needs to perform the Viterbi computation, and it therefore reduces the computation time significantly.

We view our chain as a Mealy machine [DEKM98], *i.e.* events and their emissions correspond to transitions, and paths through the chain correspond to histories. To get a one-to-one correspondence between histories and paths through our chain, we use Lemma 3 from Appendix B.1. Checking the condition of the lemma for a new event  $e_{i+1}$  involves determining whether  $(e_i, e_{i+1})$  is pairwise legal and whether condition  $\Delta(i)$  holds, both of which can be done if we know the state (not Markov state!)  $S_{i+1}$  and the position of the root of  $e_i^{\dagger}$ . So, as a first approximation to our Markov states, we can consider the set

$$\{(S, r) \mid S \text{ a legal state, and } r \in S \},$$

with the transition rule that from state  $(S, v)$ , an event  $e$  is *not* allowed only if (1)  $r(e^t) \notin S$  or (2)  $r(e^t) < v$  and  $r(e^t)$  is neither a descendant nor an ancestor of  $v$  (*i.e.*  $r(e^t)$  is *indirectly related* to  $v$ ). Now if we set

$$S'_{i+1} := S_{i+1} \setminus \{n \mid n \text{ indirectly related to } r(e_i^t) \text{ and } n < r(e_i^t)\},$$

the rule simplifies to:  $e_{i+1}$  is allowed if  $r(e_{i+1}^t) \in S'_{i+1}$ . As it stands, this definition of Markov state seemingly implies that we still need to store  $S_i$  to be able to update it to  $S_{i+1}$ , which in turn would be used to define  $S'_{i+1}$ . However, it can be shown that the straightforward update rule

$$\begin{aligned} S'_{i+1} := & \left( S'_i \setminus (e_i^t \cup \{n \in S'_i \mid n \text{ indirectly related to } r(e_i^t) \text{ and } n < r(e_i^t)\}) \right) \\ & \cup \{n \in e_i^t \mid e_i^l(n) \neq -\} \end{aligned}$$

correctly updates  $S'_i$  to  $S'_{i+1}$ , and therefore we do not need to store  $S_i$ .

We set the probability of a transition to be  $p(e)$ , the probability of the event. Note that with this choice of state, there may be several events that form transitions between a given pair of states. We can associate a graph to the current discussion, introducing a node for each state  $S'$  and joining two nodes by a set of edges, with each edge corresponding to an allowed event between the two involved states. To complete the graph, we introduce two additional nodes, one for an initial state  $A$  and the other for a final state  $\Omega$ . Furthermore, we introduce a single edge (with  $p = \prod_{n \in T} I_n$ ) between the initial state  $A$  and the state  $T$ , and join each state except  $A$  to  $\Omega$  by a single edge (with  $p = 1$ ). In this way, the probability of any path from  $A$  to  $\Omega$  through the graph is just the probability of the associated sequence of events, or history, under the TKF91 model. Moreover, since by our choice of states, each permissible history is counted exactly once, the probabilities sum up to 1. Although the resulting graph is *not* a Markov chain, since outgoing transition probabilities do not sum to 1, the above two properties are enough for the Viterbi algorithm to work.

In many situations it is useful to be able to sample from the *a posteriori* distribution of evolutionary histories. In the two-sequence case, it is a simple matter to recover the state probabilities from the collapsed probabilities  $P(i, j)$ , for example the probability for state  $S^2$  at position  $(i, j)$  in Figure 1 is  $p_2 P(i, j - 1) - p'_0 p_2 P(i - 1, j - 1)$ . This can be extended to the general case, though the details are a bit involved. A Markov chain can be used for sampling more straightforwardly, but the graph defined above for the Viterbi algorithm cannot be used as it stands. It can be turned into a Markov chain by a scalar base change, however. Let  $P_S$  be the probability of being in state  $S$ , *i.e.* the sum of probabilities of all paths ending in  $S$ . Then, the state equation can be written as  $P_S = \sum_{S'} P_{S'} t_{S'S}$ , where  $t_{S'S}$  is the transition matrix. If we introduce new variables  $\hat{P}_S = f_S P_S$ , then we can write  $\hat{P}_S = \sum_{S'} \hat{P}_{S'} \hat{t}_{S'S}$ , where  $\hat{t}_{S'S} = t_{S'S} f_S / f_{S'}$ . The graph becomes a Markov chain if  $\sum_S \hat{t}_{S'S} = 1$ . A tedious but straightforward calculation shows that  $f_S = (\prod_{n \in S} I_n)^{-1}$  solves this equation, and the transition probabilities become

$$S_i \xrightarrow{e} S_{i+1} : \frac{f_{S_{i+1}}}{f_{S_i}} p(e),$$

where  $f_A = f_\Omega = 1$ .

## 10 Discussion

The importance of incorporating statistical analysis into biological studies has become abundantly clear over the years. In particular, many interesting problems which arise in the field of bioinformatics have been successfully addressed using statistical models. For instance, the evolution of biological sequences, which is the focus of the present paper, has been formulated in a statistical framework in which mutation events are seen as stochastic processes. A closer examination of the past progress reveals, however, that although substitution processes have been modelled as continuous-time stochastic processes for more than three decades [JC69] and have been widely used in phylogenetics and geneology [Fel01], insertion and deletion processes have received less attention. A dominant reason for the observed trend is that insertion and deletion processes have only recently been generalised to arbitrary number of sequences [SH01, Hei01, Mik02, HJP02] and that the implementations of such generalisations have hitherto required quite a large running time.

In the present paper, we have constructed several algorithmic improvements which make multiple statistical alignment computationally tractable. More precisely, we have proved that the *one-state recursion* exists for arbitrary number of sequences, thus reducing the space complexity by a factor of  $O(\sqrt{5}^k)$ , where  $k$  is the number of sequences. Furthermore, we have developed a reverse traversal algorithm which calculates, with time complexity linear in the number of sequences, each transition factor appearing in the one-state dynamic programming algorithm. This reduces the time complexity of the entire algorithm to  $O(k 2^k L^k)$ , where  $L$  is the geometric mean sequence length. In fact, most of the transition factors appearing in the recursion share a common part, which hence needs to be calculated only once. Therefore, with a more careful implementation, the running time of the final algorithm can be reduced to  $O(2^k L^k)$ .

The achieved improvements are not only theoretical, but make multiple statistical alignment on trees possible in practice, as we have shown in the globin example. Compared to the ordinary hidden Markov model implementation on a four-sequence tree, our method gives rise to a speedup of about a factor 50, as well as reducing the number of states from 45 to 1, with an associated reduction in memory usage.

Our algorithm is useful mainly for trees with small number of leaves, and may be used for hypothesis testing or as a basis for the quartet method. Even for small trees, however, corner cutting techniques are necessary in practical applications. The example of §5 requires computing a table of size  $4.8 \times 10^8$  in a straightforward implementation, while corner cutting can reduce this by a factor of about 500. Because the probability mass is highly concentrated around the alignment region, we lose almost no contribution to the total probability, even though only 0.2% of the entire table is actually visited. For more than, say, 6 sequences, our method will cease to be practical, and MCMC methods will be

more useful. There is no reliable sampling method developed to date, though there were some promising attempts [HB01, JH02]. The introduced algorithm can be used to compare exact likelihood calculations and approximations, and on this basis, it might be useful at developing new sampling techniques.

The underlying sequence evolution model has the drawback of treating only single nucleotide insertions and deletions, as was already mentioned in the original paper [TKF91]. A model that can handle insertion and deletion events of whole subsequences is expected to perform better, just as affine gaps improve alignments in score-based methods. Similarly, models that incorporate some structure information will probably result in better multiple alignments as well. We hope to incorporate these ideas in the future.

## Acknowledgments

The authors wish to thank Rune Lyngsø and Alexei Drummond for helpful discussions. This research is supported by EPSRC under grant HAMJW and by MRC under grant HAMKA. Y.S. Song is supported in part by a grant from the Danish Natural Science Foundation (SNF-5503-13370).

## A States for the HJP recursion

For a fixed number  $k$  of sequences, the total number of Markov states<sup>2</sup> in the HJP recursion depends on the tree topology. Let  $T$  be a  $k$ -leaved rooted binary tree. Let  $w(n)$  denote the number of leaves joined to the internal node  $n \in T$  and let  $q(T)$  be the number of internal nodes  $n \in T$  with  $w(n) = 2$ .

If  $q(T) = 2$  and the root<sup>3</sup>  $r$  of the tree has  $w(r) = 2$ , then the number of Markov states in the HJP recursion for such a tree is equal to

$$\left( 2^k + \sum_{j=0}^{k-1} 2^j \right) + 2 \left[ \sum_{n=0}^{k-4} \left( 2^{k-2-n} + \sum_{j=0}^{k-2-n-1} 2^j \right) \right] + \sum_{n=1}^{k-4} \sum_{m=0}^{n-1} \left( 2^{n-m} + \sum_{j=0}^{n-m-1} 2^j \right).$$

As the tree topology we just considered has the least number of states, the above formula is a lower bound on the number of Markov states for any  $k$ -leaved rooted binary tree of arbitrary topology.

For  $q(T) \gg 2$ , the total number of states chiefly depends on whether or not the nodes  $n$  with  $w(n) = 2$  are assigned “#”, which stands for having a nucleotide. The number of

---

<sup>2</sup>We refer the reader to [HJP02] for their definition of “state.”

<sup>3</sup>In [HJP02], the root of a tree  $T$  is one of the internal nodes  $n \in T$  of degree 3.

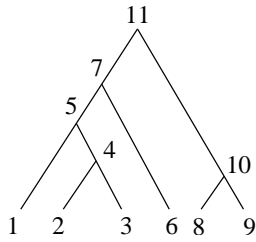


Figure 8: A 6-leaved tree with its nodes labeled according to our convention.

states associated to such choices is

$$\begin{aligned} \sum_{m=0}^q \left[ \binom{q}{m} \left( 2^{k-2m} + \sum_{j=0}^{k-2m-1} 2^j \right) \right] &= 2^k \sum_{m=0}^q \binom{q}{m} \left( 2^{-2m} + \frac{2^{k-2m} - 2}{2^k} \right) \\ &\approx 2 \cdot 2^k \left( \frac{5}{4} \right)^q. \end{aligned} \quad (5)$$

For  $q(T) = \lfloor k/2 \rfloor$ , i.e. in the case of a “balanced” tree, the expression in (5) is of order  $O(\sqrt{5}^k)$ .

## B Proofs

The main goal of this section is to prove Lemma 1 from §7.2. In the first two subsections, we first establish some results which are used in the proof of Lemma 1.

### B.1 Unique legal sequences

To define a unique ordering of events for legal histories, we introduce a condition on sequences of events. We suppose the nodes of  $T$  to be ordered. For purposes later on, we choose the specific ordering which satisfies, for every subtree  $t$  of  $T$ ,

$$\{\text{left subtree of } t\} < \{\text{right subtree of } t\} < \text{root of } t,$$

where left and right are arbitrary but fixed. An example of our choice of ordering is shown in Figure 8. We also order (context,node) pairs as follows:

$$(c, n) < (c', n') \text{ if } c = I \text{ and } c' = F, \text{ or } c = c' \text{ and } n < n'.$$

To specify a unique ordering of events, we introduce condition  $\Delta$  for a sequence of events  $E = (e_1, \dots, e_N)$  as follows:

$$\begin{aligned} \Delta(i) : \quad & t_i \cap t_{i+1} \neq \emptyset \text{ or } (c_i, r_i) < (c_{i+1}, r_{i+1}), \\ \Delta : \quad & \Delta(i) \text{ holds, } \forall 1 \leq i < N. \end{aligned}$$

**Lemma 2** *Condition  $\Delta$  holds for precisely one representative of each legal history.*

**Proof:** First we prove existence. Let  $k$  be the first index for which  $\Delta(k)$  is false. Interchanging  $e_k$  and  $e_{k+1}$  is an allowed permutation, and makes  $\Delta(k)$  true, but might render  $\Delta(k-1)$  false. In this case, interchanging  $e_{k-1}$  and  $e_k$  is allowed, and renders  $\Delta(k-1)$  and  $\Delta(k)$  true, the latter because  $\Delta(k-1)$  was true originally. Continuing backwards, eventually all  $\Delta(i)$  for  $i \leq k$  are true. By induction on  $k$  we find an allowed permutation  $\pi$  such that  $\pi(E)$  satisfies  $\Delta$ .

We now prove uniqueness. Assume there are two representatives for which condition  $\Delta$  holds,  $E$  and  $E'$ . Let  $\pi$  be the allowed permutation such that  $\pi(E) = E'$ . We can find an  $i < j$  such that  $\pi(j) + 1 = \pi(i)$  [Chr96]. Since  $\pi$  is an allowed permutation,  $t_i \cap t_j = \emptyset$ .  $E'$  satisfies  $\Delta$ , therefore  $t_j$  and  $t_i$  are part of the left and right subtree of a node  $C_0$ , respectively. Moreover,  $i + 1 \neq j$ , because this would contradict  $\Delta(i)$  for  $E$ .

We prove by induction that for each  $k \geq 0$ , a node  $C_k$  exists such that  $t_{j-k}$  is in the left subtree of  $C_k$ ,  $t_i$  is in the right subtree of  $C_k$ , and  $t_j$  is contained in either. This means that there is no  $k$  such that  $i + 1 = j - k$ , therefore there is only one representative that satisfies the  $\Delta$  condition.

We already proved that the induction hypothesis holds for  $k = 0$ . Assume it holds for certain  $k$ . Four possible cases are to be considered.

1.  $t_{j-k-1} \cap t_i \neq \emptyset$  and  $t_{j-k-1} \cap t_{j-k} \neq \emptyset$ . This implies that  $C_k \in t_{j-k-1}$  and hence  $t_{j-k-1} \cap t_j \neq \emptyset$ . Since  $\pi$  is an allowed permutation,  $\pi(i) < \pi(j-k-1)$  and  $\pi(j-k-1) < \pi(j)$  but this is impossible.

2.  $t_{j-k-1} \cap t_i \neq \emptyset$  and  $t_{j-k-1} \cap t_{j-k} = \emptyset$ . This means that  $t_{j-k-1}$  is on the right subtree of  $C_k$ , and therefore  $r_{j-k-1} > r_{j-k}$ , but this contradicts the  $\Delta$  condition, since all  $c_j = F$  because the history is legal.

3.  $t_{j-k-1} \cap t_i = \emptyset$  and  $t_{j-k-1} \cap t_{j-k} \neq \emptyset$ . This means that  $t_{j-k-1}$  is on the left part of  $C_k$ . For  $C_{k+1} = C_k$  the condition of the induction holds.

4.  $t_{j-k-1} \cap t_i = \emptyset$  and  $t_{j-k-1} \cap t_{j-k} = \emptyset$ . In this case  $r_{j-k-1} < r_{j-k} < r_i$ . The first inequality comes from the  $\Delta$  condition, while the second one is from the induction hypothesis. However, these inequalities guarantee the existence of a proper  $C_{k+1}$ . ■

In properly ordered sequences, legal sequences may be recognised by looking at pairs of events.

**Definition 2** A pair of events  $(e_1, e_2)$  is called pairwise legal if  $t_1 \supseteq t_2 \Rightarrow l_1(r_2) \neq -$ .

**Lemma 3** The set

$$\{(e_1, \dots, e_N) \mid \Delta \text{ holds, } (e_i, e_{i+1}) \text{ pairwise legal for } i = 1, \dots, N-1\}$$

consists of precisely one representative of each legal history.

**Proof:**  $\subseteq$ : Since  $\Delta$  holds, at most one of each legal history is included. We prove that illegal histories are not included. Let  $e_1, \dots, e_N$  be an illegal sequence in the set. Let  $e_i$  be

the first event with  $c_i = I$ , that is,  $r_i \notin S_i$ . Since  $c_j = F$  for  $j < i$ , in particular  $S_{i-1}$  is legal. Suppose  $t_{i-1} \cap t_i \neq \emptyset$ . If  $t_{i-1} \subset t_i$ , then  $r_i \notin S_i$  implies  $r_i \notin S_{i-1}$  but  $r_{i-1} \in S_{i-1}$  and  $r_{i-1}$  is a descendant of  $r_i$ , in contradiction with  $S_{i-1}$  legal. So  $t_{i-1} \supseteq t_i$ , but then  $r_i \notin S_i$  implies  $(e_{i-1}, e_i)$  is not pairwise legal. So  $t_{i-1} \cap t_i = \emptyset$ . Note that  $c_{i-1} = F$  and  $c_i = I$ , so  $\Delta$  is false, which is the required contradiction.

$\supseteq$ : For all legal histories there exists a sequence obeying  $\Delta$  by Lemma 2, and all neighbouring pairs in legal sequences are pairwise legal.  $\blacksquare$

This result enables us to determine the sequences of events we want to sum over, by looking only at neighbouring pairs. To use the inclusion/exclusion trick we next have to characterise sequences of event pairs that do *not* obey our condition.

## B.2 Disallowed sequences and illegal sets

**Definition 3** Suppose  $S_i$  is a legal state, and suppose that for all  $j = i, \dots, k-1$  we have  $(e_j, e_{j+1})$  not pairwise legal or  $\Delta(j)$  false. We call such (sub)sequences disallowed sequences.

**Lemma 4** Let  $(e_i, e_{i+1}, \dots, e_k)$  be a disallowed sequence starting from state  $S_i$ . Then:

- a) The sequence  $(c_j, r_j)$ ,  $i \leq j \leq k$ , is strictly decreasing.
- b) All events  $e_j$  with  $c_j = F$  are disjoint.
- c) For each node  $n$ , there is at most one event  $e_j$  with  $n \in t_j$  and  $l_j(n) \neq -$ .

Note that statement (c) implies that there is at most one emission per leaf.

**Proof:** If  $(e_j, e_{j+1})$  is not pairwise legal, then  $r_{j+1}$  is a descendant of  $r_j$  and  $l_j(r_{j+1}) = -$ . This means that  $r_j > r_{j+1}$  and  $c_{j+1} = I$ , so  $(c_j, r_j) > (c_{j+1}, r_{j+1})$ . If  $\Delta(j)$  is false, then it directly follows from the definition of condition  $\Delta$  that  $(c_j, r_j) > (c_{j+1}, r_{j+1})$ .

Now, that the sequence  $(c_j, r_j)$  is strictly decreasing implies that, for some  $m$ , we have  $c_i = c_{i+1} = \dots = c_m = F$  and  $c_{m+1} = c_{m+2} = \dots = c_k = I$ . For  $i \leq j < m$  the pair  $(e_j, e_{j+1})$  is pairwise legal since  $c_{j+1} = F$ , therefore we must have  $\Delta(j)$  false, which implies that  $t_j \cap t_{j+1} = \emptyset$ . Furthermore,  $t_j \cap t_{j+1} = \emptyset$  and  $r_j > r_{j+1}$  together imply that there exists a subtree  $t$  of  $T$  such that  $t_j$  is in a right subtree of  $t$ , whereas  $t_{j+1}$  is in a left subtree of  $t$ . Applying this reasoning sequentially, we conclude that for all  $i \leq a < b \leq m$  there exists a subtree  $t$  such that  $t_a$  ( $t_b$ ) is in a right (left) subtree of  $t$ . Hence,  $t_a \cap t_b = \emptyset$  and it follows that all of  $t_i, t_{i+1}, \dots, t_m$  are pairwise disjoint.

Note that  $S_{m+1}$  is a legal state since  $S_i$  is a legal state and all events  $e_i, e_{i+1}, \dots, e_m$  are legal events. For every  $m < j \leq k$ , that the event  $e_j$  is illegal implies that  $r_j \notin S_j$ . Now, suppose  $S_j \cap t_j \neq \emptyset$ . Then,  $S_j$  must contain at least one descendant of  $r_j$ , which is possible only if, since  $S_{m+1}$  is a legal state and  $r_j \notin S_j$ , at least one of  $e_{m+1}, e_{m+2}, \dots, e_{j-1}$  has a

birth at a descendant node of  $r_j$ . But, this contradicts the fact that  $r_{m+1} > r_{m+2} > \dots > r_k$ , and therefore we must have  $S_j \cap t_j = \emptyset$  for every  $m < j \leq k$ .

Finally, the results from the previous two paragraphs together imply (c).  $\blacksquare$

**Lemma 5** *For a given legal initial state  $S$ , there is a one-to-one correspondence between illegal sets and disallowed sequences starting from  $S$ .*

**Proof:** We prove the correspondence by establishing two injections.

From disallowed sequence to illegal set: The events with  $c_j = F$  are pairwise disjoint, and their trees not the subtree of any other because the initial state is legal. The trees with  $c_j = I$  may be subtree of another  $t_k$ , in which case  $l_k(r_j) = -$ , so the events form an illegal set. There is only one disallowed sequence giving rise to this set, since the sequence  $(c_j, r_j)$  corresponding to the disallowed sequence is ordered.

From illegal set to disallowed sequence: Take the events whose trees are not subset of another, and whose root is in  $S$ , sorted decreasing by root. Follow them by the other events, sorted decreasing by root. A disallowed sequence results.  $\blacksquare$

### B.3 Proof of Lemma 1

For the ease of reference, we here restate the lemma. We denote by  $E(\mathbf{K})$  a set of legal sequences of events which emit the prefixes of the given sequences  $A^1, \dots, A^k$  up to position  $\mathbf{K} = (K_1, \dots, K_k)$ , and which contains precisely one representative of each history. The quantity of interest is the likelihood

$$P(\mathbf{K}) := \sum_{(e_1, \dots, e_N) \in E(\mathbf{K})} p(e_1) \cdots p(e_N).$$

**Lemma 1** *The likelihood  $P(\mathbf{K})$  satisfies the following equation:*

$$\begin{aligned} P(\mathbf{K}) &= \sum_e P(\mathbf{K} - \mathbf{v}_e) p(e) \\ &- \sum_{\{e_1, e_2\} \text{ illegal}} P(\mathbf{K} - \mathbf{v}_{e_1} - \mathbf{v}_{e_2}) p(e_1) p(e_2) \\ &+ \sum_{\{e_1, e_2, e_3\} \text{ illegal}} P(\mathbf{K} - \mathbf{v}_{e_1} - \mathbf{v}_{e_2} - \mathbf{v}_{e_3}) p(e_1) p(e_2) p(e_3) \\ &- \dots \end{aligned} \tag{6}$$

Here  $\mathbf{v}_e$  is the emission vector corresponding to the event  $e$ , that is, a  $k$ -dimensional vector whose component corresponding to each leaf is 0 if the leaf is labeled  $-$ , or 1 otherwise.

**Proof:** Let  $S(e)$  denote the state after event  $e$  in state  $S$ . Define  $P_S(\mathbf{K})$  to be the likelihood of emitting legal sequences up to position  $\mathbf{K}$  and ending up in state  $S$ . We allow illegal states  $S$ , in which case  $P_S(K) = 0$  always. For brevity, we denote by  $C(S_0, e_0, S_1, e_1, \dots, e_{n-1}, S_n)$



the condition  $S_0(e_0) = S_1 \wedge \dots \wedge S_{n-1}(e_{n-1}) = S_n$ , and by  $\text{Dis}(S, e_0, \dots, e_n)$  the condition that  $(e_0, \dots, e_n)$  is a disallowed sequence when starting from state  $S$ . Then

$$\begin{aligned}
P_S(\mathbf{K}) &= \sum_{(S',e):C(S',e,S)} P_{S'}(\mathbf{K} - \mathbf{v}_e)p(e) \\
&- \sum_{\substack{(S',e,S'',e'): \\ C(S'',e',S',e,S) \wedge \text{Dis}(S'',e',e)}} P_{S''}(\mathbf{K} - \mathbf{v}_e - \mathbf{v}_{e'})p(e')p(e) \\
&+ \sum_{\substack{(S',e,S'',e',S''',e''): \\ C(S''',e'',S'',e',S',e,S) \wedge \\ \wedge \text{Dis}(S''',e'',e',e)}} P_{S'''}(\mathbf{K} - \mathbf{v}_e - \mathbf{v}_{e'} - \mathbf{v}_{e''})p(e'')p(e')p(e) \\
&- \dots
\end{aligned}$$

This recursion sums over all sequences of events which do not have disallowed subsequences, since all terms in the first line that include a disallowed subsequence (of length 2) are subtracted in the second line; this also subtracts terms involving disallowed subsequences of length 3 that were not included in the first place, which are added in again in the third line, and so on. By Lemma 3, this implies that we indeed sum over all sequences of legal events (which end up in state  $S$ ). Note that for illegal states  $S$ ,  $P_S(\mathbf{K}) = 0$  always, so we may define  $\text{Dis}(S, e_0, \dots)$  arbitrarily for illegal states  $S$ .

Now we sum the recursion equation over all states  $S$ . The left-hand side becomes  $P(\mathbf{K})$ . The first summation on the right-hand side turns into

$$\sum_S \sum_{(S',e):S'(e)=S} P_{S'}(\mathbf{K} - \mathbf{v}_e)p(e) = \sum_{(S',e)} P_{S'}(\mathbf{K} - \mathbf{v}_e)p(e) = \sum_e P(\mathbf{K} - \mathbf{v}_e)p(e),$$

which is the desired term. Similarly, the second summation becomes

$$\sum_{(S'',e',e):\text{Dis}(S'',e',e)} P_{S''}(\mathbf{K} - \mathbf{v}_e - \mathbf{v}_{e'})p(e')p(e).$$

Observe that the summand is independent of the ordering of events  $e', e$ . Hence, we can also sum over illegal *sets* of events and over  $S''$ , using Lemma 5. This turns the summation into

$$\sum_{S''} \sum_{\{e',e\} \text{ illegal}} P_{S''}(\mathbf{K} - \mathbf{v}_e - \mathbf{v}_{e'})p(e')p(e) = \sum_{\{e',e\} \text{ illegal}} P(\mathbf{K} - \mathbf{v}_e - \mathbf{v}_{e'})p(e')p(e),$$

which is the desired term. All other summations are dealt with analogously. Note that only finitely many summands contribute to the recursion, since illegal sets have a bounded number of elements. ■

## References

- [Chr96] D. A. Christie. Sorting permutations by block-interchanges. *Information Processing Letters*, 60:165–169, 1996.

- [DEKM98] R. Durbin, S. Eddy, A. Krogh, and G. Mitchison. *Biological sequence analysis*. Cambridge University Press, 1998.
- [Edd01] S. Eddy. Hmmer: Profile hidden markov models for biological sequence analysis (<http://hammer.wustl.edu/>), 2001.
- [Fel81] J. Felsenstein. Evolutionary trees from dna sequences: a maximum likelihood approach. *J. Mol. Evol.*, 17:368–376, 1981.
- [Fel01] J. Felsenstein. Phylip: phylogeny software package. University of Washington, 2001.
- [Gol98] N. Goldman. Effects of sequence alignment procedures on estimates of phylogeny. *BioEssays*, 20:287–290, 1998.
- [HB01] I. Holmes and W. J. Bruno. Evolutionary HMMs: a Bayesian approach to multiple alignment. *Bioinformatics*, 17(9):803–820, 2001.
- [Hei01] J. Hein. An algorithm for statistical alignment of sequences related by a binary tree. In R. B. Altman, A. K. Dunker, L. Hunter, K. Lauderdale, and T. E. Klein, editors, *Pac. Symp. Biocomp.*, pages 179–190. World Scientific, 2001.
- [Hir77] D. S. Hirshberg. Algorithms for the longest common subsequence problem. *J. ACM*, 24:664–675, 1977.
- [HJP02] J. Hein, J. L. Jensen, and C. N. S. Pedersen. Recursions for statistical multiple alignment. Technical Report 425, Dept. of Theor. Stat., Univ. of Aarhus, January 2002.
- [HWK<sup>+</sup>00] J. Hein, C. Wiuf, B. Knudsen, M. B. Møller, and G. Wibling. Statistical alignment: Computational properties, homology testing and goodness-of-fit. *J. Mol. Biol.*, 302:265–279, 2000.
- [JC69] T.H. Jukes and C.R. Cantor. *Mammalian Protein Metabolism*, chapter Evolution of protein molecules, pages 21–132. Academic Press, 1969.
- [JH02] J.L. Jensen and J. Hein. Gibbs sampler for statistical multiple alignment. Technical Report 429, Dept. of Theor. Stat., Univ. of Aarhus, September 2002.
- [Lee01] M. S. Y. Lee. Unalignable sequences and molecular evolution. *Trends in Ecology & Evolution*, 16:681–685, 2001.
- [Mik02] I. Miklós. An improved algorithm for statistical alignment of sequences related by a star tree. *Bul. Math. Biol.*, 64:771–779, 2002.
- [NW70] S. B. Needleman and C. D. Wuntch. A general method applicable to the search for similarities in the amino acid sequences in two proteins. *J. Mol. Biol.*, 48:443–453, 1970.

- [SH01] M. Steel and J. Hein. Applying the Thorne-Kishino-Felsenstein model to sequence evolution on a star-shaped tree. *Appl. Math. Let.*, 14:679–684, 2001.
- [SK83] D. Sankoff and J. B. Kruskall. *Time warps, string edits and macromolecules: the theory and practice of sequence comparison*. Addison-Wesley, 1983.
- [TKF91] J. L. Thorne, H. Kishino, and J. Felsenstein. An evolutionary model for maximum likelihood alignment of DNA sequences. *J. Mol. Evol.*, 33:114–124, 1991.
- [WLG01] S. Whelan, P. Lió, and N. Goldman. Molecular phylogenetics: state-of-the-art methods for looking into the past. *Trends in Genetics*, 17:262–272, 2001.