

A simple method for finding related sequences by adding probabilities of alternative alignments

Martin C. Frith^{1,2,3}

¹Artificial Intelligence Research Center, AIST, Tokyo, Japan

²Department of Computational Biology and Medical Sciences, University of Tokyo, Chiba, Japan

³Computational Bio Big Data Open Innovation Laboratory, AIST, Tokyo, Japan

Abstract

The main way of analyzing genetic sequences is by finding sequence regions that are related to each other. There are many methods to do that, usually based on this idea: find an alignment of two sequence regions, which would be unlikely to exist between unrelated sequences. Unfortunately, it is hard to tell if an alignment is likely to exist by chance. Also, the precise alignment of related regions is uncertain. One alignment does not hold all evidence that they are related. We should consider alternative alignments too. This is rarely done, because we lack a simple and fast method that fits easily into practical sequence-search software.

Here is described a simplest-conceivable change to standard sequence alignment, which sums probabilities of alternative alignments. Remarkably, this makes it easier to tell if a similarity is likely to occur by chance. This approach is better than standard alignment at finding distant relationships, at least in a few tests. It can be used in practical sequence-search software, with minimal increase in implementation difficulty or run time. It generalizes to different kinds of alignment, e.g. DNA-versus-protein with frameshifts. Thus, it can widely contribute to finding subtle relationships between sequences.

Introduction

Many methods have been developed for finding related parts of nucleotide or protein sequences. They

usually follow a standard approach: define positive and negative scores for letter matches, mismatches, and gaps, then seek high-scoring alignments. Why is this bizarre, intricate approach a good way to find related sequences? Its result depends utterly on the score parameters, so how should we choose them?

The answer to both questions is that the approach is equivalent to probability-based alignment (Dayhoff et al. 1978; Durbin et al. 1998). This uses probabilities of matches, mismatches, and gaps, e.g. their average rates in the kinds of related sequences we wish to find, and seeks high-probability alignments. An alignment's score is a log probability ratio (Frith 2020):

$$\text{score} = \log \frac{\text{prob}(\text{alignment})}{\text{prob}(\text{null alignment})}. \quad (1)$$

$\text{prob}(\text{alignment})$ is the product of probabilities of the gaps, aligned letters, and unaligned letters in both sequences. A “null alignment” is a length-zero alignment between the two sequences. A precise definition of $\text{prob}(\text{null alignment})$ was given previously (Frith 2020), and is repeated in this paper's Supplemental Methods. A high alignment score indicates high probability with these match/mismatch/gap rates. The base of the log (e.g. \log_2 , \log_{10}) is unspecified, because any base is equally correct. A change of base merely rescales the scores. (In the precursor to this paper (Frith 2020), $\log x$ was equivalently written as $t \ln x$, where t is an arbitrary positive constant.) This approach is termed “pairwise local alignment”, where “local” means that the alignment is not necessarily of

the whole sequences.

A more direct way to judge whether two sequence regions are related is to calculate their probability without fixing an alignment. In other words, use a score like this:

$$\begin{aligned} \text{score} &= \log \frac{\sum_{\text{alignments}} \text{prob}(\text{alignment})}{\text{prob}(\text{null alignment})} \\ &= \log \sum_{\text{alignments}} \frac{\text{prob}(\text{alignment})}{\text{prob}(\text{null alignment})}. \end{aligned} \quad (2)$$

If the sum includes all possible ways of aligning the sequence regions, it gives their total probability with the assumed rates of matches, mismatches, and gaps. Thus, it includes all evidence that the sequence regions are related, while one alignment includes just some evidence.

A similar sum-over-alignments method has been used with great success in the HMMER software (Eddy 2008). So it is curious that many sequence-search tools developed since HMMER do not use this approach. No-one has described a sum-over-alignments method that fits easily into typical sequence-search software, with minimal increase in implementation difficulty or computational cost. HMMER is too complicated and specialized, e.g. it has an enormous “implicit probabilistic model”, and sequence-length-dependent parameters (Eddy 2008). HMMER is designed to search a smallish “query” sequence against a large “database”, so it can’t find related parts of two large sequences, e.g. two genomes.

It is important to know whether a similarity score (eqs. 1, 2) is likely to occur by chance between unrelated sequences. The simplest definition is: the probability of an equal or higher score between two sequences of random i.i.d. (independent and identically distributed) letters. Even this simplest definition, however, is hard to calculate. There is a solution for gapless alignment, in the limit where the sequences are long (Karlin and Altschul 1990). Gapped alignment relies on conjecture and simulation. BLAST can calculate it only for a few sets of score parameters, not including realistic DNA scores that distinguish transition (a:g, c:t) from transversion mismatches (Altschul et al. 1997). Some widely-used aligners do not calculate it at all (Harris 2007). The ALP library

can calculate it for arbitrary parameters, by complex simulations (Sheetlin et al. 2016). Remarkably, HMMER calculates it in an easier way, based on other conjectures. The scope of those conjectures is unclear, e.g. they seem to hold only after tweaking HMMER to use a “uniform entry/exit distribution” (Eddy 2008).

Here is described a minimal modification of standard alignment, which sums probabilities of alternative alignments (eq. 2). It can replace the alignment component of practical sequence-search software, with minimal increase in implementation difficulty or run time. This also makes it easier to calculate probabilities of similarity scores between random sequences, based on a clear conjecture.

Limitations of previous methods

Many previous studies have summed probabilities of alternative alignments. This section surveys some representatives, to understand why this approach is rarely used to find related parts of sequences.

One line of research considered substitutions, insertions, and deletions occurring over time, so that match/mismatch/gap probabilities depend on the length of time (Bishop and Thompson 1986; Thorne et al. 1991; Thorne and Churchill 1995; Hein et al. 2000; Knudsen and Miyamoto 2003; Rivas and Eddy 2015). These methods achieve several aims. Given two related sequences, they can: infer substitution/insertion/deletion rates per unit time, find the most likely alignment, and report the probability that any pair of letters descend from a common ancestor. Furthermore, Hein et al. used a log likelihood ratio to compare two hypotheses, that two whole sequences are related or not:

$$\log \frac{\text{prob}(\text{sequences} \mid \text{related})}{\text{prob}(\text{sequences} \mid \text{unrelated})}. \quad (3)$$

The numerator is a sum of probabilities of all possible alignments between the whole sequences.

Another line of research (including the present study) simplifies things by ignoring time-dependence, and directly using probabilities of matches, mismatches, and gaps. Allison et al. (1992) estimated such probabilities from a pair of related sequences, allowing complex gap-length distributions. They also

used a log likelihood ratio to compare hypotheses that two whole sequences are related or not.

In a pioneering study, Bucher and Hofmann (1996) used a log likelihood ratio to compare two hypotheses, that two sequences have a related segment or not:

$$\text{score} = \log \frac{\text{prob}(\text{sequences} \mid \text{have related part})}{\text{prob}(\text{sequences} \mid \text{lack related part})} \quad (4)$$

Each hypothesis has probabilities for the lengths of the two sequences. Their two hypotheses are identical in this regard. Thus, the sequence lengths never favor either hypothesis. This is arguably unnatural, biologically and mathematically. Biologically, if related regions are equally likely to occur per unit sequence length, longer sequences are more likely to have related regions. For example, a longer chromosome is more likely to contain a mobile DNA element. Mathematically, they need two summation algorithms. The likelihood ratio becomes a fraction where the numerator is a sum of exponentiated alignment scores, and the denominator is another sum. This denominator increases with increasing sequence lengths (fig. 2 in Webb et al. 2002). So scores of related regions decrease with increasing lengths of sequences containing them.

HMMER uses a log likelihood ratio similar to eq. 4 (Eddy 2008). HMMER is more general, however. It finds related parts between a sequence and a sequence family: the family has position-specific letter and gap probabilities. The null hypothesis (no related parts) is compared to either a “unihit” hypothesis (one related part), or a “multihit” hypothesis (one or more related parts). Each hypothesis has probabilities for the length of the sequence. The hypotheses are roughly equal in this regard. HMMER achieves that by adjusting each hypothesis to fit the length of each sequence that is provided for analysis (Eddy 2008). Again, scores of related regions decrease with increasing sequence length. HMMER’s DNA version, nhmmer, applies this analysis to sequence windows around potential matches in long sequences (Wheeler and Eddy 2013).

Lunter et al. (2008) compared DNA sequences accurately, using probability methods similar to the present study. In their scenario, they are given a pair of sequences that are assumed to have one related

region. So they didn’t use any score to find related regions.

FEAST is practical heuristic software for finding related sequence parts, and it sums probabilities of alternative alignments (Hudek and Brown 2010). It doesn’t determine whether similarities would be unlikely to exist between random sequences. It is slower than standard alignment, and implementation difficulty was not made clear.

Several of these previous methods use a likelihood ratio, which is the optimal way to judge between two hypotheses with those likelihoods (Neyman and Pearson 1933). It is however not clear that comparing two hypotheses is the best way to find related sequence parts. For example, if we compare the human and mouse X chromosomes, it is not very useful just to know whether or not they have 0 related parts. There is a tradeoff between judging whether sequences are related, and how they are related. A single alignment specifies a relationship in detail, but has the lowest power to detect subtle relationships. Summing over all alignments has the highest power, but least detail. We seek a compromise between these extremes.

The final previous idea that we shall examine is “hybrid alignment” (Yu and Hwa 2001). When hybrid alignment compares a sequence of length m to a sequence of length n , it considers $m \times n$ hypotheses for how they are related. Each hypothesis concerns the the length- i prefix of one sequence and length- j prefix of the other. The hypothesis is that the sequences have a related region that ends exactly at the ends of these prefixes. The probability associated with this hypothesis is the sum of probabilities of all alignments that start anywhere and end at the ends of the prefixes. This sum of alignment probabilities is converted to a score (eq. 2). Remarkably, it is easy to calculate the probability of an equal or higher score between random i.i.d. sequences, provided that the alignment parameters satisfy a “conservation condition”. This calculation was not proven correct, but was supported by theoretical arguments and tests on random sequences.

The present study builds on hybrid alignment. Hybrid alignment has been neglected, perhaps because its description was mathematically intricate, and it was presented as a “semi-probabilistic”, “hybrid” method.

Algorithm 1 Optimum alignment extension in an $n \times n$ block (lower-right gray block in Fig. 1A)

$$\begin{array}{llll} X_{00} \leftarrow 0 & X_{i0} \leftarrow -\infty & X_{0i} \leftarrow -\infty & (0 < i \leq n) \\ v \leftarrow -\infty & Z_{i0} \leftarrow -\infty & Y_{0i} \leftarrow -\infty & (0 \leq i \leq n) \end{array}$$

$$\left. \begin{array}{l} w \leftarrow \max(X_{ij}, Y_{ij}, Z_{ij}) \\ v \leftarrow \max(v, w) \\ X_{i+1\ j+1} \leftarrow w + S_{R_{i+1}Q_{j+1}} \\ Y_{i+1\ j} \leftarrow \max(w + a_D, Y_{ij} + b_D) \\ Z_{i\ j+1} \leftarrow \max(w + a_I, Z_{ij} + b_I) \end{array} \right\} \begin{array}{l} 0 \leq i \leq n \\ 0 \leq j \leq n \end{array}$$

(X_{ij}), or R_i aligned to a gap (Y_{ij}), or Q_j aligned to a gap (Z_{ij}). Thus, w is the highest possible score for a global (end-to-end) alignment of R_1, \dots, R_i and Q_1, \dots, Q_j . v tracks the best alignment score seen so far, and is the final output. Several variants of this algorithm are possible. This variant has simple boundary conditions, while minimizing add, max, and array access operations (Cameron et al. 2004). The alignment score after left and right extension is:

$$\text{score} = v_{\text{left}} + (\text{core score}) + v_{\text{right}}. \quad (5)$$

New extension algorithm

This section shows the new algorithm (alg. 2), and the next section explains how it sums probabilities. There are two changes from alg. 1. The first, which makes no difference to the result, is to work with exponentiated values: probability ratios instead of log probability ratios. This is so we can easily add them. This change means that $-\infty \Rightarrow 0$, $0 \Rightarrow 1$, and addition \Rightarrow multiplication. The other change is: maximization \Rightarrow addition. This calculates the sum of probability ratios for all alignment extensions in the gray area (eq. 2), instead of the maximum probability ratio. The score with left and right extension is:

$$\text{score} = \log[v'_{\text{left}}] + (\text{core score}) + \log[v'_{\text{right}}]. \quad (6)$$

The core score is the same as in eq. 5. The core is typically a short, high-similarity, gapless alignment.

Algorithm 2 Sum over alignment extensions in an $n \times n$ block (lower-right gray block in Fig. 1A)

$$\begin{array}{llll} X'_{00} \leftarrow 1 & X'_{i0} \leftarrow 0 & X'_{0i} \leftarrow 0 & (0 < i \leq n) \\ v' \leftarrow 0 & Z'_{i0} \leftarrow 0 & Y'_{0i} \leftarrow 0 & (0 \leq i \leq n) \end{array}$$

$$\left. \begin{array}{l} w' \leftarrow X'_{ij} + Y'_{ij} + Z'_{ij} \\ v' \leftarrow v' + w' \\ X'_{i+1\ j+1} \leftarrow w' \cdot S'_{R_{i+1}Q_{j+1}} \\ Y'_{i+1\ j} \leftarrow w' \cdot a'_D + Y'_{ij} \cdot b'_D \\ Z'_{i\ j+1} \leftarrow w' \cdot a'_I + Z'_{ij} \cdot b'_I \end{array} \right\} \begin{array}{l} 0 \leq i \leq n \\ 0 \leq j \leq n \end{array}$$

Alternative alignments of the core region are not considered.

Note that the number of computational operations is unchanged: the only changes are replacing addition with multiplication, and maximization with addition. So the computational complexity is unchanged. The practical run time depends on implementer expertise (e.g. using SIMD), but plausibly it need not increase too much.

Unfortunately, the values (probability ratios) in alg. 2 may become too large for the computer to handle. This can be fixed by occasionally rescaling them (Supplemental Methods). Because the rescaling is only done occasionally, it does not increase the run time noticeably.

Parameters from probabilities

To use alg. 1 or 2, we must choose values for the parameters (e.g. a'_D , b'_I). They are related to probabilities of matches, mismatches, and gaps. We need a precise definition of these probabilities. We shall use the definition in Fig. 2. The main alignment probabilities are shown in the middle: α_D is the deletion start probability, β_D the deletion extension probability, α_I the insertion start probability, and β_I the insertion extension probability. Apart from starting an insertion or deletion, the other possibilities are to have a pair of aligned letters (probability γ), or end the alignment (probability $1 - \gamma - \alpha_D - \alpha_I$). A pair of

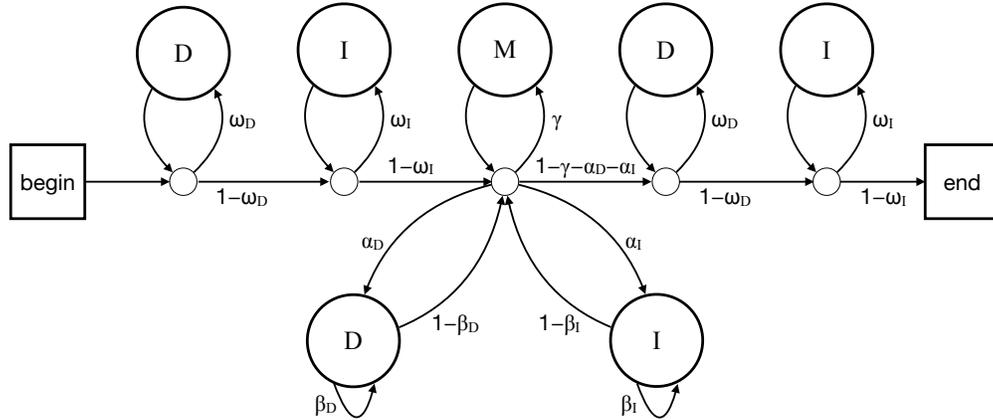


Figure 2: A probability model for two sequences with a related region. The model produces different alignments with different probabilities. Starting at “begin”, the arrows are followed according to their probabilities (e.g. ω_D versus $1 - \omega_D$). Each pass through a **D** state produces an unaligned letter x in one sequence (R), with probabilities ϕ_x . Each pass through an **I** state produces an unaligned letter y in the other sequence (Q), with probabilities ψ_y . Each pass through **M** generates two aligned letters, x in R and y in Q , with probabilities π_{xy} .

aligned letters has probability π_{xy} of being letter type x in sequence R and y in sequence Q . There are also probabilities for letters left and right of the alignment: ω_D and ω_I (for sequence R and Q respectively). An unaligned letter is of type z with probability ϕ_z in R , and ψ_z in Q .

A “null alignment” is produced by a path that never traverses the α_D , α_I , or γ arrows. This corresponds to a length-0 alignment between the sequences. The score of a null alignment is traditionally zero, which is achieved by eq. 1. This means that an alignment score is determined solely by the aligned regions, and not by flanking sequences or their lengths.

Having defined these probabilities, we can state the parameters for alg. 2, as shown previously (Frith 2020):

$$S'_{xy} = \frac{\gamma}{\omega_D \omega_I} \cdot \frac{\pi_{xy}}{\phi_x \psi_y} \quad (7)$$

$$a'_D = \alpha_D(1 - \beta_D)/\omega_D \quad (8)$$

$$b'_D = \beta_D/\omega_D \quad (9)$$

$$a'_I = \alpha_I(1 - \beta_I)/\omega_I \quad (10)$$

$$b'_I = \beta_I/\omega_I \quad (11)$$

With these parameters, alg. 2 calculates the sum of

probability ratios in eq. 2. The parameters for alg. 1 are defined similarly (Supplemental Methods).

Fig. 2 is not the only way to define the probabilities (Frith 2020). There are other ways that cause no change to alg. 1, but cause slight changes to alg. 2. The definition in Fig. 2 seems to produce the simplest equations, and minimizes the number of computational operations in alg. 2. This is a minor improvement over hybrid alignment, which defined probabilities with an unnatural asymmetry between insertions and deletions (fig. S2 in Frith 2020).

A non-heuristic similarity score

What would we do if we did not need fast heuristics? As mentioned in the Introduction, we seek a compromise between using just one alignment, and summing over all alignments. We shall use a compromise with two useful features: it is approximated by the heuristic approach (Fig. 1B), and we can tell if a similarity is unlikely to exist between random sequences.

When we compare a length- m sequence (R_1, R_2, \dots, R_m) to a length- n sequence (Q_1, Q_2, \dots, Q_n), we shall consider $(m+1) \times (n+1)$ hypotheses for how they are related. Imagine cutting R into a prefix of

length i and suffix of length $m-i$, and Q into a prefix of length j and suffix of length $(n-j)$. These cuts define a point, illustrated in Fig. 1C by the touching corners of the gray rectangles. The hypothesis is that the sequences are related by any alignment that passes through this point (including alignments that start or end there). The score of this hypothesis is given by eq. 2, with the sum taken over these alignments.

We can calculate this score for each (i, j) . We first calculate W_{ij}^F , the sum of probability ratios of all alignments that end at the ends of the prefixes:

$$\begin{aligned} 0 \leq i \leq m, \quad 0 \leq j \leq n \\ W_{ij}^F &= S'_{R_i Q_j} W_{i-1, j-1}^F + Y_{i-1, j}^F + Z_{i, j-1}^F + 1 \\ Y_{ij}^F &= a'_D W_{ij}^F + b'_D Y_{i-1, j}^F \\ Z_{ij}^F &= a'_I W_{ij}^F + b'_I Z_{i, j-1}^F \end{aligned}$$

The boundary condition is: when $i < 0$ or $j < 0$, $W_{ij}^F = Y_{ij}^F = Z_{ij}^F = 0$.

We then calculate W_{ij}^B , the sum of probability ratios of all alignments that start at the starts of the suffixes:

$$\begin{aligned} m \geq i \geq 0, \quad n \geq j \geq 0 \\ W_{ij}^B &= S_{R_{i+1} Q_{j+1}} W_{i+1, j+1}^B + Y_{i+1, j}^B + Z_{i, j+1}^B + 1 \\ Y_{ij}^B &= a'_D W_{ij}^B + b'_D Y_{i+1, j}^B \\ Z_{ij}^B &= a'_I W_{ij}^B + b'_I Z_{i, j+1}^B \end{aligned}$$

The boundary condition is: when $i > m$ or $j > n$, $W_{ij}^B = Y_{ij}^B = Z_{ij}^B = 0$. Finally, the sum of probability ratios of all alignments passing through (i, j) is $W_{ij}^F W_{ij}^B$. So the score is $\log[W_{ij}^F W_{ij}^B]$.

A minor point is that this score definition excludes alignments that pass through (i, j) while traversing the β_D or β_I arrow (Fig. 2). This makes it closer to typical heuristic methods that extend alignments from a high-quality gapless “core” (Fig. 1).

Conservation condition

We wish to know the probabilities of similarity scores between random sequences. Because our scores, $\log[W_{ij}^F W_{ij}^B]$, are similar to those of hybrid alignment, we shall assume (and test) that we can use the same

approach. This approach requires that the alignment parameters satisfy a “conservation condition”.

The conservation condition is related to alignment length bias. The probabilities in Fig. 2 may be biased towards short or long alignments between two given sequences. For example, if α_D , α_I , and γ are all low (almost 0), but ω_D and ω_I are high (almost 1), shorter alignments will have higher probability. In the converse situation, longer alignments are favored. It was shown previously (Frith 2020) that the probabilities are unbiased when:

$$\frac{\gamma}{\omega_D \omega_I} + \frac{\alpha_D(1 - \beta_D)}{\omega_D - \beta_D} + \frac{\alpha_I(1 - \beta_I)}{\omega_I - \beta_I} = 1. \quad (12)$$

Next, let us see the conservation condition. Suppose we compare two random i.i.d. sequences, with letter probabilities Φ_x and Ψ_y . The conservation condition is

$$\left(\sum_{x,y} \Phi_x \Psi_y S'_{xy} \right) + \frac{a'_D}{1 - b'_D} + \frac{a'_I}{1 - b'_I} = 1. \quad (13)$$

In practice, we shall always assume that $\Phi_x = \phi_x$ and $\Psi_y = \psi_y$. In that case, eq. 13 is the same as eq. 12.

The conservation condition depends on the precise definition of gap probabilities (Fig. 2), as explained previously (Frith 2020). The conservation condition for hybrid alignment was different (eq. 28 in Yu and Hwa 2001), because they defined gap probabilities differently (fig. S2 in Frith 2020). The conservation condition generalizes an equation in the proven solution for probabilities of gapless alignment scores (eq. 4 in Karlin and Altschul 1990).

Similarity scores occurring by chance

We can calculate the probability of a similarity score s between two random i.i.d. sequences: one with length m and letter probabilities ϕ_x , the other with length n and probabilities ψ_y . Chance similarities occur at a constant rate between any part of one sequence and any part of the other. The p -value is the probability of a score $\geq s$ occurring, and the E -value is the expected number of distinct similarities with score $\geq s$. In the limit where m and n are large, the number

of distinct similarities follows a Poisson distribution, which means that $p = 1 - e^{-E}$. The prediction is that

$$E = Kmn / \exp(s), \quad (14)$$

and therefore

$$p = 1 - e^{-Kmn / \exp(s)}. \quad (15)$$

Here, \exp is defined to mean inverse of \log .

These formulas have one unknown parameter, K . We can estimate K by generating some pairs of random i.i.d. sequences, calculating $s_{\max} = \log[\max_{ij}(W_{ij}^F W_{ij}^B)]$ for each pair, and fitting K (Supplemental Methods).

Implementation details

This approach to finding related sequence parts was added to the LAST software. This implementation is just a proof of principle: the heuristics are not necessarily the best. LAST can use either alg. 1 with E -values from the ALP library, or alg. 2. In both cases, it gets a representative alignment by running alg. 1 and tracing back a way to get the maximum v . LAST runs these algorithms not in an $n \times n$ block (Fig. 1A), but in a range that is adjusted as the algorithm proceeds (Fig. 1B, Supplemental Methods).

Sometimes, LAST finds overlapping alignments from two different “cores” (Fig. 1B). To remove such redundancy: if two representative alignments share an endpoint (same coordinates in both sequences), LAST keeps just one, with highest similarity score.

LAST estimates K from random sequences. It is not clear how long these sequences need to be. Based on some tests (Supplemental Fig. S1), LAST uses 50 pairs of length-500 sequences by default.

The match/mismatch/gap probabilities (π_{xy} , α_D , etc.) were determined using LAST-TRAIN, which finds related regions of given sequences and infers these probabilities (Hamada et al. 2017). LAST-TRAIN was modified to make the probabilities satisfy eq. 12. It sets ω_D and ω_I by assuming that they are equal, and finding the unique value that satisfies eq. 12 (Supplemental Methods). LAST-TRAIN is not necessarily the best way to get the probabilities. As an alternative, a LAST-TRAIN option was added, which makes it work

as usual except that the letter probabilities (π , ϕ , and ψ) are fixed to those of a BLOSUM or PAM matrix.

Results

Similarity scores occurring by chance

Let us test whether we can calculate probabilities of similarity scores between random sequences. Three sequence-search scenarios are considered, with typical DNA, at-rich DNA, and proteins.

The first scenario is finding ancient (shared with reptiles) repeats in the human genome. This can be done by comparing the genome to repeat consensus sequences from Dfam (Storer et al. 2021). Some are relics of transposable elements, others are source unknown. The match/mismatch/gap probabilities were found by comparing the genome to the consensus sequences with LAST-TRAIN. Then, pairs of random i.i.d. sequences were generated, with letter probabilities ϕ_x and ψ_y equal to those found by LAST-TRAIN. For each pair of sequences, the highest score s_{\max} was found by the non-heuristic algorithm. The observed scores are consistent with their probabilities predicted by eq. 15 (Fig. 3 left).

The next scenario was to find related parts of the *Plasmodium falciparum* and *Plasmodium yoelii* genomes, which are both about 80% a+t. Alignment parameters were found by LAST-TRAIN, then similarity scores were found between random at-rich sequences (Fig. 3 middle). Again, the observed scores are consistent with eq. 15.

The third scenario was to find related parts of proteins from *Aquifex aeolicus* (a heat-loving bacterium) and *Pyrolobus fumarii* (a heat-loving archaeon). Alignment parameters were found by LAST-TRAIN, then similarity scores were found between random sequences (Fig. 3 right). Yet again, the observed scores are consistent with eq. 15. This remains true if we use BLOSUM62 letter probabilities (Supplemental Fig. S3A).

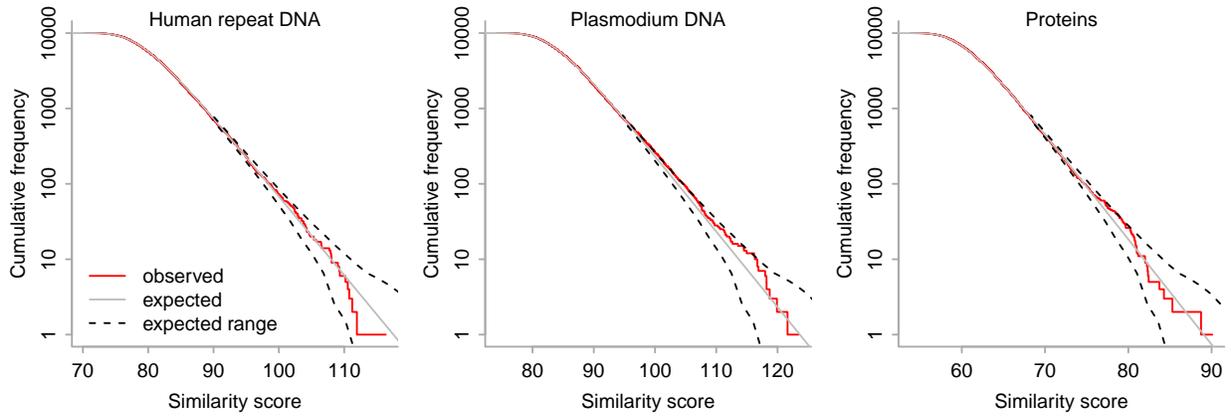


Figure 3: Non-heuristic similarity score (s_{max}) between 10000 pairs of random i.i.d. sequences (length 10000), for three sets of alignment parameters. The frequency has a 5% chance of being outside the dashed lines (2.5% each).

```

A  tccttcctctcttttctcccttcctctctctcccttcacttccttttctcttctccctc      C  acttgatcttagccaaaaggccgagaagcgat
   | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
   tntttcccccctctcctaccttntctctctctccattggcttttcttttctctgcccctcttc
                                     acttgatcttagccaaaaggccgagaagcgat

B  cctgtttggctcctggagtcttagcagagaccaccagcaggggagctaattactgcccag-----tctgtaatgtggacctccgtcc
   || | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
   ccattcagtccttggcctctctgctgagactgccaacaagggtgccaattagtccaattgtgatggtggttttgtgatgtggacacctgtcc

```

Figure 4: Alignment between parts of: **A** reversed human Chromosome Y (upper), and UCON1 (lower). This is the highest-scoring alignment in the top-left panel of Fig. 5. **B** Human Chromosome 6 (upper) and LFSINE_Vert (lower). **C** Human Chromosome 21 (upper) and human U2 (lower).

Related versus biased sequences

A high similarity score could occur because the sequence regions have a common ancestor, and/or similar composition bias (Fig. 4A). We can investigate this by comparing two biological sequences after reversing (but not complementing) one of them (but see Glidden-Handgis and Wheeler 2024). Then there are no segments related by evolution, but there may be similarities due to composition bias (Fig. 4A).

Thus, the reversed human genome was searched against the repeat consensus sequences, the reversed *P. yoelii* genome against *P. falciparum*, and reversed *P. fumarii* proteins against *A. aeolicus* proteins. In each case, alg. 2 found more high similarity scores than expected between random sequences (Fig. 5 upper row), due to biased regions (Fig. 4A, Supplemental Fig. S4).

For standard alignment, these similarities can be avoided by detecting biased regions with `tantan` (Frith 2011a), and “masking” them. An effective masking scheme is to change $S_{xy} \leftarrow \min(S_{xy}, 0)$ when x or y is in a biased region (Frith 2011b). For alg. 2, this masking scheme was tried: $S'_{xy} \leftarrow \min(S'_{xy}, 1)$.¹ With this masking scheme, the similarity scores between reversed and non-reversed sequences are close to what is expected between random sequences (Fig. 5 lower row).

The new method can be more sensitive

Related sequence parts were sought by standard alignment (alg. 1) and by the new method (alg. 2), with biased-sequence masking. The two methods often find identical alignments, but with different scores and E -values. The new method tended to assign lower E -values, for alignments of the human genome to repeat consensus sequences (Fig. 6A), and *P. fumarii* proteins to *A. aeolicus* proteins (Fig. 6B, Supplemental Fig. S3B). For example, both methods found the same alignment between human Chromosome 6 and the LFSINE_Vert consensus (Fig. 4B): the alg. 1 E -value is 22, and the alg. 2 E -value is 0.000031. The alg. 2 score is 185.3, much more than the highest score with reversed sequences (147.5, Fig. 5 lower left). This

¹Another untested idea is $S'_{xy} \leftarrow \min(S'_{xy}, \frac{\gamma}{\omega_D \omega_I})$.

means the new method can find related sequence parts more confidently.

In other words, at a fixed E -value cutoff, the new method predicts more related regions. The preceding results indicate that the E -values are not overconfident. This implies that the new method finds more true positives for a given false positive rate.

As a further test, similar regions were sought after shuffling the letters of the sequences used for Fig. 6A. As expected, similarities were found with E -value ≥ 1 , and the alg. 2 E -values did not trend lower than the alg. 1 E -values (Supplemental Fig. S5).

The new method can be less sensitive

Another sequence search scenario was considered: finding U2 fragments in the human genome. U2 is one of the small nuclear RNAs that splice introns. The genome has many fragmentary copies of it, which are thus considered repeats.

The human genome was searched against a U2 DNA sequence, using the match/mismatch/gap probabilities found earlier for the human genome and ancient repeat consensus sequences. This time, alg. 2 tended to assign higher E -values than alg. 1 (Fig. 6C). For example, both methods found the alignment shown in Fig. 4C: the alg. 1 E -value is 0.00021, and the alg. 2 E -value is 0.0011.

The likely explanation is as follows. The genome has many short, near-exact copies of the first 30–40 bases of U2. The match/mismatch/gap probabilities, however, reflect greater sequence divergence (e.g. Fig. 4B). These mistuned probabilities pessimize both algorithms, but they pessimize alg. 2 more. That is because probabilities with greater divergence attach more weight to alternative alignments.

This suggests that alg. 2 should work better than alg. 1 if the probabilities are a good fit to the sequences. To test this, match/mismatch/gap probabilities were found by comparing the genome to the U2 sequence with LAST-TRAIN. Then, the genome was searched against U2 using these probabilities. This time, alg. 2 tends to assign slightly lower E -values than alg. 1 (Fig. 6D). Even with well-tuned probabilities, the advantage of alg. 2 is expected to decrease for more

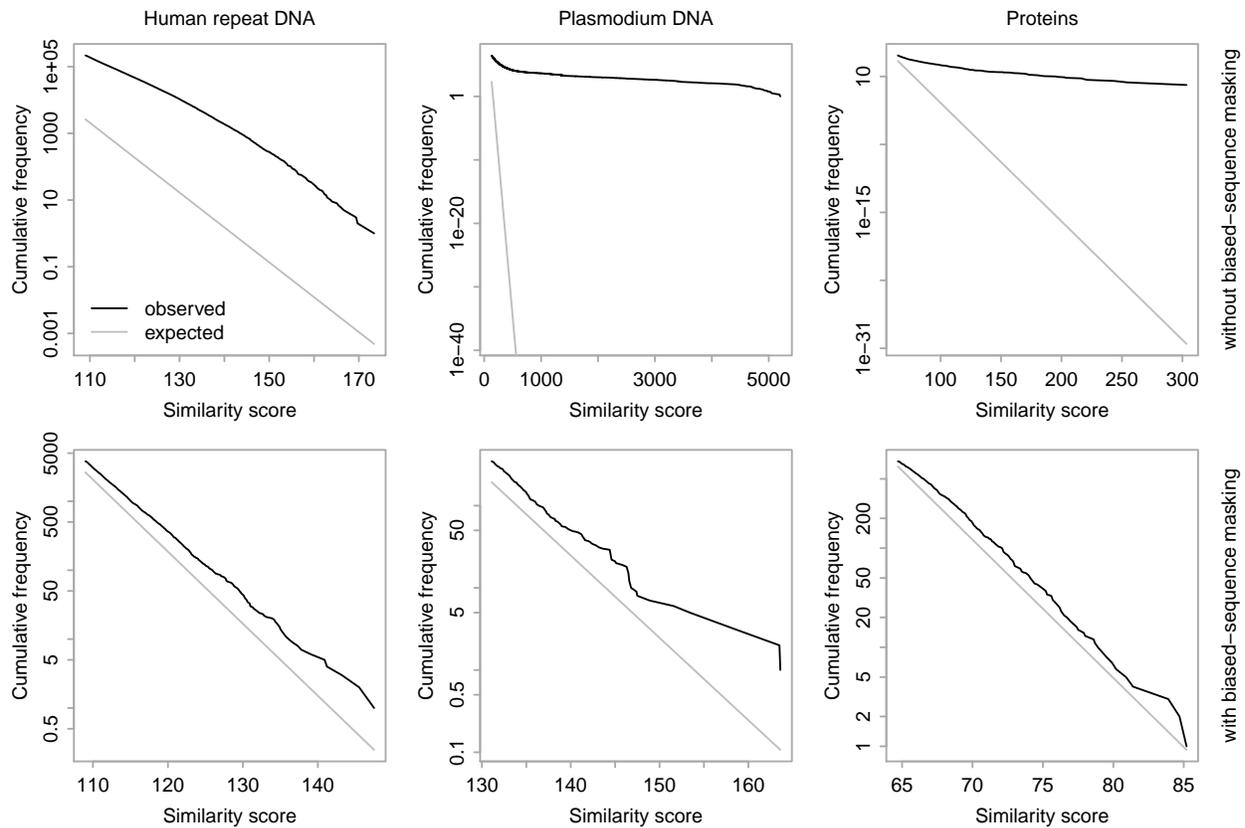


Figure 5: Heuristic similarity scores, between reversed and non-reversed sequences, with or without biased-sequence masking.

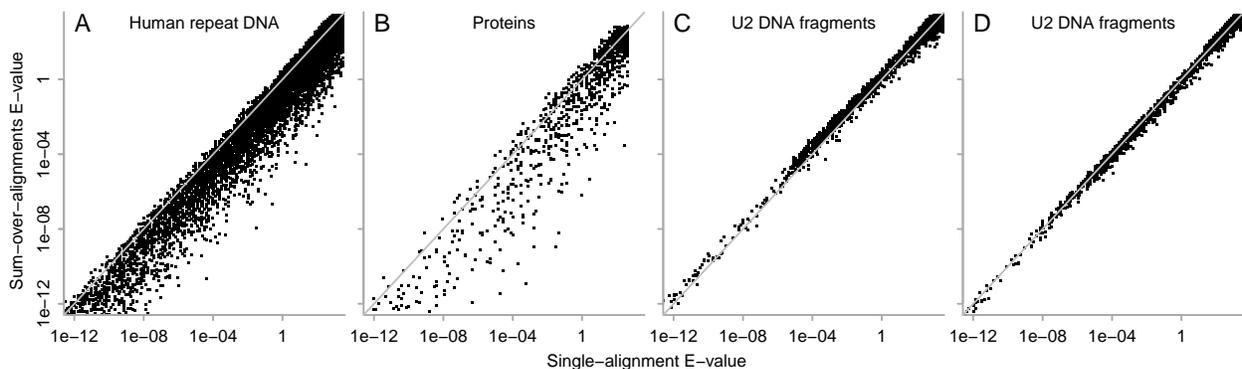


Figure 6: E -values for identical alignments found by alg. 1 (horizontal axis), and alg. 2 (vertical axis). Each point is one alignment. The diagonal gray lines indicate equal E -values. These E -values were calculated with m = total length of all reference sequences (e.g. all *A. aeolicus* proteins), and n = total length of all query sequences (e.g. all *P. fumarii* proteins). For DNA, one query strand was searched against both reference strands, so m was multiplied by 2.



Figure 7: An alignment (the same as Fig. 4B) with colors indicating the probability that each column is correct. ~ indicates the “core” (Fig. 1), whose column probabilities are assumed to be 1. * indicates non-core exact matches.

similar sequences, because the evidence for relatedness is more concentrated in one alignment.

Discussion

We have seen how to find related parts of sequences, by summing probabilities of alternative alignments between them. This is better than standard alignment at detecting subtle relationships, at least in a few tests. It can be worse than standard alignment, when the probabilities are tuned for high divergence but the sequences have low divergence. The probabilities should be tuned for different use cases, e.g. using LAST-TRAIN. The new method is more beneficial for higher divergence, where alternative alignments contribute more. In a spectacular example, it found gene-regulating DNA conserved in diverse animals since the Precambrian (Frith and Ni 2023).

A major further benefit is that the new method simplifies E -value calculation. This makes it easier to experiment with different alignment parameters, beyond the BLAST parameter sets. This can be useful for proteins (Ng et al. 2000), but especially for DNA, because there has been less effort to optimize DNA parameters.

Moreover, the new method can be generalized to other kinds of alignment. For example, we used the same approach for DNA-versus-protein alignment (Yao and Frith 2023). We defined probabilities similarly to Fig. 2, but with extra probabilities for frameshifts. We summed probabilities of alternative alignments, and accurately predicted similarity scores between random sequences. This method was highly sensitive: it found ancient and degraded relics of Paleozoic mobile elements in vertebrate genomes. Some previous methods have used two kinds of alignment gap: short-and-frequent and long-and-rare (Allison

et al. 1992; Lunter et al. 2008). This can be done by elaborating Fig. 2. The approach of this paper is predicted to work in that case too. This makes it easier to try such different kinds of alignment, in particular because it provides E -values. Another alignment elaboration is position-specific match, mismatch, and gap rates (Durbin et al. 1998; Eddy 2008; Steinegger et al. 2019; Roddy et al. 2024): this was incorporated into hybrid alignment (Yu et al. 2002).

By combining alg. 2 with another similar algorithm (Supplemental Methods), we can calculate the probability that each alignment column is correct (Fig. 7), based on the assumed rates of matches, mismatches, and gaps (Allison et al. 1992; Durbin et al. 1998). It is also possible to make alignments based on these probabilities: for example, align letters whose alignment probability is > 0.5 (Miyazawa 1995).

For short sequences, the E -values are too high. Eq. 14 assumes the number of chance similarities is proportional to mn . That becomes too high for short sequences, because a similarity must have non-zero length and fit in the sequences. There are corrections for this effect, including in hybrid alignment, but it is somewhat complex (Yu and Hwa 2001).

We should consider some basic aims and assumptions. Suppose we seek related regions between a genome of length m , and q “query” sequences with lengths n_1, n_2, \dots, n_q . We could assume that related regions are equally likely to occur per query sequence, or per unit sequence length. We could aim to find relationships for as many query sequences as possible, or as many related regions as possible. If part of query i is similar to part of the genome, we can calculate an E -value from eq. 14 with sequence lengths m and n_i . This increasingly disfavors related regions in increasingly long queries. As mentioned in the Introduction, some previous methods similarly penalize related regions in longer sequences. The alternative is

to use eq. 14 with sequence lengths m and $\sum_{j=1}^q n_j$. This treats related regions equally, regardless of the containing sequence's length. The aim is important when testing sequence search methods, because the test may count the number of sequences or related regions found.

Software availability

LAST is available at <https://gitlab.com/mcfrith/last>, and as Supplemental Code. The LAST-TRAIN output files from this study are at <https://gitlab.com/mcfrith/sum-align>. There too is a simple Python script to help use these ideas. It takes as input α_D , α_I , β_D , β_I , and γ , finds $\omega_D = \omega_I$ that satisfies eq. 12, and outputs a'_D , a'_I , b'_D , b'_I , and $\frac{\gamma}{\omega_D \omega_I}$. LAST includes code for calculating the non-heuristic similarity score (used for estimating K), in a self-contained file (`mcf_alignment_path_adder.cc`).

Competing interest statement

The author declares no competing interests.

Acknowledgments

I am grateful to John Spouge for encouragement that the E -values would be simple, Yi-Kuo Yu for advice about hybrid alignment, and Travis Wheeler for advice about HMMER. This research was supported by the Japan Science and Technology Agency [JP-MJCR21N6].

References

- Allison L, Wallace C, Yee C. 1992. Finite-state models in the alignment of macromolecules. *Journal of Molecular Evolution* **35**: 77–89.
- Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ. 1997. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Research* **25**: 3389–3402.
- Bishop M, Thompson EA. 1986. Maximum likelihood alignment of DNA sequences. *Journal of molecular biology* **190**: 159–165.
- Bucher P, Hofmann K. 1996. A sequence similarity search algorithm based on a probabilistic interpretation of an alignment scoring system. In: *Proc Int Conf Intell Syst Mol Biol*. Pp. 44–51.
- Cameron M, Williams HE, Cannane A. 2004. Improved gapped alignment in BLAST. *IEEE/ACM Transactions on Computational Biology and Bioinformatics* **1**: 116–129.
- Dayhoff M, Schwartz R, Orcutt B. 1978. A model of evolutionary change in proteins. *Atlas of protein sequence and structure* **5**: 345–352.
- Durbin R, Eddy SR, Krogh A, Mitchison G. 1998. *Biological sequence analysis: probabilistic models of proteins and nucleic acids*. Cambridge University Press.
- Eddy SR. 2008. A probabilistic model of local sequence alignment that simplifies statistical significance estimation. *PLoS Computational Biology* **4**: e1000069.
- Frith MC. 2011a. A new repeat-masking method enables specific detection of homologous sequences. *Nucleic Acids Research* **39**: e23–e23.
- Frith MC. 2011b. Gentle masking of low-complexity sequences improves homology search. *PLoS One* **6**: e28819.
- Frith MC. 2020. How sequence alignment scores correspond to probability models. *Bioinformatics* **36**: 408–415.
- Frith MC, Ni S. 2023. DNA conserved in diverse animals since the Precambrian controls genes for embryonic development. *Molecular Biology and Evolution*: msad275.
- Glidden-Handgis G, Wheeler TJ. 2024. WAS IT A MATch I SAW? Approximate palindromes lead to overstated false match rates in benchmarks using reversed sequences. *Bioinformatics Advances* **4**: vbae052.
- Guidi G, Ellis M, Rokhsar D, Yelick K, Buluç A. 2021. BELLA: Berkeley efficient long-read to long-read aligner and overlapper. In: *SIAM Conference on Applied and Computational Discrete Algorithms (ACDA21)*. SIAM: pp. 123–134.
- Hamada M, Ono Y, Asai K, Frith MC. 2017. Training alignment parameters for arbitrary sequencers with LAST-TRAIN. *Bioinformatics* **33**: 926–928.
- Harris RS. 2007. “Improved pairwise alignment of genomic DNA.” PhD thesis. The Pennsylvania State University.
- Hein J, Wiuf C, Knudsen B, Møller M, Wibling G. 2000. Statistical alignment: computational properties, homology testing and goodness-of-fit. *Journal of Molecular Biology* **302**: 265–279.

- Hudek AK, Brown DG. 2010. FEAST: sensitive local alignment with multiple rates of evolution. *IEEE/ACM Transactions on Computational Biology and Bioinformatics* **8**: 698–709.
- Karlin S, Altschul SF. 1990. Methods for assessing the statistical significance of molecular sequence features by using general scoring schemes. *Proceedings of the National Academy of Sciences* **87**: 2264–2268.
- Knudsen B, Miyamoto MM. 2003. Sequence alignments and pair hidden Markov models using evolutionary history. *Journal of molecular biology* **333**: 453–460.
- Liu D, Steinegger M. 2023. Block Aligner: an adaptive SIMD-accelerated aligner for sequences and position-specific scoring matrices. *Bioinformatics*: btad487.
- Lunter G, Rocco A, Mimouni N, Heger A, Caldeira A, Hein J. 2008. Uncertainty in homology inferences: assessing and improving genomic sequence alignment. *Genome research* **18**: 298–309.
- Ma B, Tromp J, Li M. 2002. PatternHunter: faster and more sensitive homology search. *Bioinformatics* **18**: 440–445.
- Miyazawa S. 1995. A reliable sequence alignment method based on probabilities of residue correspondences. *Protein Engineering, Design and Selection* **8**: 999–1009.
- Neyman J, Pearson ES. 1933. On the problem of the most efficient tests of statistical hypotheses. *Philosophical Transactions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character* **231**: 289–337.
- Ng PC, Henikoff JG, Henikoff S. 2000. PHAT: a transmembrane-specific substitution matrix. *Bioinformatics* **16**: 760–766.
- Noé L, Kucherov G. 2004. Improved hit criteria for DNA local alignment. *BMC Bioinformatics* **5**: 1–9.
- Rivas E, Eddy SR. 2015. Parameterizing sequence alignment with an explicit evolutionary model. *BMC Bioinformatics* **16**: 1–23.
- Roddy JW, Rich DH, Wheeler TJ. 2024. nail: software for high-speed, high-sensitivity protein sequence annotation. *bioRxiv*. DOI: [10.1101/2024.01.27.577580](https://doi.org/10.1101/2024.01.27.577580).
- Roytberg M, Gambin A, Noé L, Lasota S, Furetova E, Szczurek E, Kucherov G. 2009. On subset seeds for protein alignment. *IEEE/ACM Transactions on Computational Biology and Bioinformatics* **6**: 483–494.
- Sheetlin S, Park Y, Frith MC, Spouge JL. 2016. ALP & FALP: C++ libraries for pairwise local alignment E-values. *Bioinformatics* **32**: 304–305.
- Steinegger M, Meier M, Mirdita M, Vöhringer H, Haunsberger SJ, Söding J. 2019. HH-suite3 for fast remote homology detection and deep protein annotation. *BMC Bioinformatics* **20**: 1–15.
- Storer J, Hubley R, Rosen J, Wheeler TJ, Smit AF. 2021. The Dfam community resource of transposable element families, sequence models, and genome annotations. *Mobile DNA* **12**: 1–14.
- Suzuki H, Kasahara M. 2017. Acceleration of nucleotide semi-global alignment with adaptive banded dynamic programming. *bioRxiv*. DOI: [10.1101/130633](https://doi.org/10.1101/130633).
- Thorne JL, Churchill GA. 1995. Estimation and reliability of molecular sequence alignments. *Biometrics*: 100–113.
- Thorne JL, Kishino H, Felsenstein J. 1991. An evolutionary model for maximum likelihood alignment of DNA sequences. *Journal of Molecular Evolution* **33**: 114–124.
- Webb BJM, Liu JS, Lawrence CE. 2002. BALSAs: Bayesian algorithm for local sequence alignment. *Nucleic Acids Research* **30**: 1268–1277.
- Wheeler TJ, Eddy SR. 2013. nhmmer: DNA homology search with profile HMMs. *Bioinformatics* **29**: 2487–2489.
- Yao Y, Frith MC. 2023. Improved DNA-versus-protein homology search for protein fossils. *IEEE/ACM Transactions on Computational Biology and Bioinformatics* **20**: 1691–1699.
- Yu YK, Bundschuh R, Hwa T. 2002. Hybrid alignment: high-performance with universal statistics. *Bioinformatics* **18**: 864–872.
- Yu YK, Hwa T. 2001. Statistical significance of probabilistic sequence alignment and related local hidden Markov models. *Journal of Computational Biology* **8**: 249–282.
- Zhang Z, Berman P, Miller W. 1998. Alignments without low-scoring regions. *Journal of Computational Biology* **5**: 197–210.



A simple method for finding related sequences by adding probabilities of alternative alignments

Martin C Frith

Genome Res. published online August 16, 2024

Access the most recent version at doi:[10.1101/gr.279464.124](https://doi.org/10.1101/gr.279464.124)

Supplemental Material <http://genome.cshlp.org/content/suppl/2024/09/12/gr.279464.124.DC1>

P<P Published online August 16, 2024 in advance of the print journal.

Accepted Manuscript Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.

Creative Commons License This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see <https://genome.cshlp.org/site/misc/terms.xhtml>). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at <http://creativecommons.org/licenses/by-nc/4.0/>.

Email Alerting Service Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or [click here](#).



To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>
