

Two Simple and Efficient Algorithms to Compute the SP-Score Objective Function of a Multiple Sequence Alignment

Vincent Ranwez

► To cite this version:

Vincent Ranwez. Two Simple and Efficient Algorithms to Compute the SP-Score Objective Function of a Multiple Sequence Alignment. PLoS ONE, 2016, 11 (8), 10.1371/journal.pone.0160043 . hal-01594194

HAL Id: hal-01594194 https://hal.science/hal-01594194

Submitted on 26 Sep 2017

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License



GOPEN ACCESS

Citation: Ranwez V (2016) Two Simple and Efficient Algorithms to Compute the SP-Score Objective Function of a Multiple Sequence Alignment. PLoS ONE 11(8): e0160043. doi:10.1371/journal. pone.0160043

Editor: Yang Zhang, University of Michigan, UNITED STATES

Received: May 30, 2016

Accepted: July 12, 2016

Published: August 9, 2016

Copyright: © 2016 Vincent Ranwez. This is an open access article distributed under the terms of the <u>Creative Commons Attribution License</u>, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: All relevant data are within the paper.

Funding: This work has been supported by the French research agency: Agence Nationale de la Recherche (ANR-10- BINF-01-02 "Ancestrome") (http://www.agence-nationale-recherche.fr/).

Competing Interests: The authors have declared that no competing interests exist.

RESEARCH ARTICLE

Two Simple and Efficient Algorithms to Compute the SP-Score Objective Function of a Multiple Sequence Alignment

Vincent Ranwez*

Montpellier SupAgro, UMR AGAP, 34060, Montpellier, France

* Vincent.ranwez@supagro.fr

Abstract

Background

Multiple sequence alignment (MSA) is a crucial step in many molecular analyses and many MSA tools have been developed. Most of them use a greedy approach to construct a first alignment that is then refined by optimizing the sum of pair score (SP-score). The SP-score estimation is thus a bottleneck for most MSA tools since it is repeatedly required and is time consuming.

Results

Given an alignment of n sequences and L sites, I introduce here optimized solutions reaching O(nL) time complexity for affine gap cost, instead of $O(n^2L)$, which are easy to implement.

Introduction

A wide range of molecular analyses rely on multiple sequence alignments (MSA), e.g., prediction of tridimensional structures [1], phylogenetic inference [2] or detection of positive selection [3]. In all these studies, the initial MSA can strongly impact conclusions and biological interpretations [4,5]. As a consequence, MSA is a richly developed area of bioinformatics and computational biology.

Most MSA software use a greedy approach to construct a first alignment that is then refined by optimizing the sum of pair score (SP-score) [$\underline{6},\underline{7}$] using a 2-cut strategy. This strategy consists in partitioning the current solution into two sub-alignments that are subsequently realigned; the resulting MSA replaces the previous one if its SP score is improved [$\underline{7},\underline{8},\underline{9},\underline{10}$]. The SP-score estimation is thus a bottleneck for most MSA tools since it is repeatedly required and is time consuming as existing solutions have a time complexity of $O(n^2L)$ for an alignment of nsequences and L sites. The practical importance of having an efficient solution to compute SPscore lies within the following 'Muscle software paper' [$\underline{8}$] quotation: "Notice that computation of the SP score dominates the time complexity of refinement and of MUSCLE overall[...]. It is natural to seek an O(nL) expression for [SP-score estimation...], but to the best of our knowledge no solution is known." The main result of this paper is an optimized algorithmic solution to estimate SP-score for affine gap costs in O(nL). I also introduce a more versatile solution able to handle more general gap cost penalty functions in $O(nL + n^2G_{max})$, with $G_{max} \leq L$ being the maximum number of gap intervals within one aligned sequence.

Materials and Methods

Definitions and notations

A multiple sequence alignment (MSA) for a set of sequences $\{s_1 \dots s_n\}$, defined with alphabet Σ , is a set of *n* sequences $\{S_1 \dots S_n\}$ which are defined on an enriched alphabet $\Sigma \cup \{`-'\}$ so that all S_i have the same length *L* and, $\forall i$, removing '-' from S_i leads to s_i . The aim of MSA tools is to position gaps (stretches of '-') so that characters at the same position (constituting a site) are (most likely) homologous. This is usually achieved through a heuristic optimization of the sum of pair score (SP-score). The SP-score of an MSA is obtained by considering all pairwise alignments it induced. Given two sequences S_i and S_j of MSA \mathcal{A} the corresponding pairwise restriction ($\mathcal{A} \mid \{S_i, S_j\}$) is the alignment made up of two sequences S'_i and S'_j obtained by removing the '-' of S_i (resp. S_j) whenever S_j (resp. S_i) also has a gap at this position/site. The score of this pairwise alignment is the sum of the substitution/match scores induced by sites without gaps, plus the sum of scores associated to the gap intervals (maximal sub-sequences of consecutive gap symbols) observed in S'_i and S'_j .

The SP-score of an alignment is classically obtained in $O(n^2L)$ by summing up the score of its $\binom{n}{2}$ induced pairwise alignments (Algorithm 1 and 2), and can be decomposed into two terms: SPs, the contribution of substitutions/matches and, SPg the contribution of induced gap intervals (denoted IG').

In molecular biology $|\Sigma|$ is a constant so typically small (4 for nucleotides and 20 for amino acids) that $L|\Sigma|^2$ and $n|\Sigma|^2$, compared to nL, can safely be ignored in asymptotic complexity analysis. Under this assumption, all solutions described here have an O(nL) space complexity.

Algorithm. 1. Basic $O(n^2L)$ computation of the SP-score of an alignment \mathcal{A} of n sites and L sequences given subst(.,.) and g_cost(.) functions. The subst(.,.) function provides, in O(1), the elementary score for two non gap characters on the same site, e.g. using BLOSUM matrix [11]. The g_cost(.) function provides, also in O(1), the cost of a gap interval based on its position and size. The compute_gap_intervals(.) subroutine (Algorithm 2) returns in O(|S|) the list of the gap intervals of its input sequence S.

<u>Algorithm1</u>: compute_SP_score

<u>Input</u>: -The *n* aligned sequences S_i of alignment \mathcal{A}

-a function subst(x,y) returning in O(1) the score for two non gap characters x and y on the same site of ${\cal A}$

```
-a function g_cost(IG') returning in O(1) the cost of a gap interval IG'
Output: the SP score of this alignment
```

$$\begin{split} & \text{SP}_{\text{s}} = 0 \text{; } \text{SP}_{\text{g}} = 0 \\ & \text{for } S_{j} \text{ in } S_{1} \dots S_{n} \\ & \text{for } S_{j} \text{ in } S_{i+1} \dots S_{n} \\ & S'_{i} = S'_{j} = "" \\ & \text{for } k \text{ in } 1 \dots L \\ & \text{ if } (\text{not } (S_{i}[k] == `-' \text{ and } S_{j}[k] == `-')) \\ & S'_{i} = S'_{i} + S_{i}[k] \\ & S'_{j} = S'_{j} + S_{j}[k] \\ & \text{ if } (S_{i}[k] \neq '-' \text{ and } S_{j}[k] \neq '-') \\ & \text{ SP}_{\text{s}} = \text{SP}_{\text{s}} + \text{ subst } (S'_{i}[k], S'_{j}[k]) \\ & \text{ for } \text{IG' in compute_gap_intervals } (S'_{i}) \cup \text{ compute_gap_intervals } (S'_{j}) \\ & \text{ SP}_{\text{g}} = \text{SP}_{\text{g}} + \text{g_cost } (\text{IG'}) // \text{ e.g., } \text{g_cost } (\text{IG'}) = \{ \text{ return } \text{gap}_{0} \\ & + \text{ IG' [length] * gap_{\text{ext}} } \end{split}$$

```
Return SP<sub>s</sub> + SP<sub>g</sub>
```

Algorithm. 2. An O(L) algorithm to compute the list of gap intervals, ordered by their gap start position, of a sequence *S* of length *L*. Note that, thought IG[length] is not explicitly set, it is assumed it can be access since *IG*[length] is simply *IG*[end]-*IG*[start]+1.

Algorithm 2: compute gap intervals

Input: An aligned sequences S_i of length L

Output: The list of gap intervals of S_i , ordered by gap start position

 $LG = \{\}$; // the list of gap intervals of S_i found so far

IG = NULL; // the current gap interval

for *k* in 1...*L*

if $(S_i[k] == '-' \text{ and } IG == \text{NULL}) // \text{ start a new gap interval}$

IG = new Interval (start = k)if $(S_i[k] \neq '-' \text{ and } IG \neq \text{NULL}) // \text{the current gap interval finish at previous}$

position

IG end] = k-1; append a copy of IG to the list LG IG = NULL

if (IG≠NULL) // handle terminal gap if any IG[end] = L; append a copy of IG to the list LG Return LG

Efficient algorithm to compute SP-score using general gap cost penalties

The SPs part of the SP-score can be computed in O(nL) by using a table of size $L|\Sigma|$ containing for each site the number of each (non gap) symbol (e.g. [8]). This strategy does not work for SPg except for the basic, but unrealistic, constant gap cost where $g_cost(IG') = IG'[length].gap$ cost. I introduce here a more efficient solution to the SP-score computation problem accounting for most gap function penalties (including affine, log, log-affine penalties). The main idea is to pre-compute the list of gap intervals of each sequence S_i , ordered by gap start position, this can easily be done in O(L) using Algorithm 2. The compute_gap_intervals(S'_i) and compute_gap_intervals(S'_i) of Algorithm 1, observed in $\mathcal{A}|\{S_i,S_i\}$, can then be efficiently deduced by processing the gaps of compute_gap_intervals(S_i) \cup compute_gap_intervals(S_i) according to order of opening (as done to merge two sorted lists in linear time, e.g. during a merge step of the 'merge sort' algorithm) while maintaining the number of gaps facing gaps encountered so far (i.e. the shift between S and S' site coordinate systems for current position). The resulting SPscore algorithm (Algorithm 3) has a complexity of $O(nL + n^2 G_{max})$, with $G_{max} \leq \lfloor L/2 \rfloor$ the maximum number of gap intervals within one aligned sequence, instead of $O(n^2L)$. Note that the difference with the naïve algorithm is especially important when sequences contain few long gap stretches but that in the worst case, where most sequences have a number of gap intervals close to L, this algorithm has the same $O(n^2L)$ complexity as the naïve solution.

Optimal algorithm to compute SP-score using affine gap cost penalties

Affine gap penalties (where $g_cost(IG') = gap_O + IG'[length].gap_{ext}$) are frequently used. For such gap penalties, the total of gap extension penalties (SP_{ge}) can also be efficiently computed in O(nL), by counting the number of gaps per site. However, gap opening cannot be counted exactly based on local site information (e.g [8]) only approximated. Though pessimistic gap count approximation [9] is often used during the dynamic programming steps producing new candidate alignments, the exact SP score is generally preferred to decide which alignments are better than the current one. Algorithm 4 provides a simple and exact solution to compute the SP-score under affine gap penalties in O(nL), which is also the time complexity for just reading an alignment of *n* sequences of *L* sites.

The key idea is to note that a gap IG_i will add a number of gap opening penalties equal to n minus the number of interval IG_j so that $IG_i \subseteq IG_j$. In order to find out how many gaps encompass IG_i , sites are processed from left to right while maintaining an array indicating, for each left site, the number of gap stretches already opened at this position and not yet closed. For all gaps IG_i closing at the current position i, the value stored at index $IG_t[leb]$ of this array provides, in O(1), the number of gap stretches encompassing IG_i ; before considering site i+1, this array is maintained updated by decreasing by 1 all values stored at indices between $IG_t[leb]$ and i for all IG_i ending at i—hence updates for all sites overall require O(nL). External and internal gaps are often penalized with different affine functions. The proposed O(nL) solution can handle this refinement by: firstly, using different characters (e.g. '-' and '_) to represent the two different gap types while computing SP_{ge}; and secondly, testing each gap interval type in Algorithm 3 (using gap interval start/end positions) to select the adequate gap opening cost.

Algorithm. 3. Given the gap interval lists LG_i , LG_j of sequences $S_i \in A$ and $S_j \in A$; this algorithm returns in $O(|LG_i| + |LG_j|)$ the restricted gap interval lists LG'_i , LG'_j that would be

```
observed in \mathcal{A}|\{S_i, S_i\} without actually building this restricted alignment.
 Algorithm 3: compute pairwise restricted gap intervals
 Input:-LG_i, LG_j the ordered lists of gap intervals for S_i \in A and S_i \in A
 <u>Output</u>:-LG'_{i}, LG'_{i} the lists of gap intervals in \mathcal{A} \mid \{S_{i}, S_{j}\}
 IG_i = first(LG_i); IG_j = first(LG_j)
 shift = 0;
 LG'_{i} = LG'_{i} = \{ \}
 IG'_{i} = \text{new Interval}(\text{start} = -1)
 IG'_i = new Interval (start = -1) // using -1 allows to check if interval start
 has already been set or not
 while (IG_i \neq NULL and IG_i \neq NULL)
   if (IG_{i}[\text{start}] == IG_{i}[\text{start}])
       if (IG_i = IG_i) / / both intervals disappear when A is restricted to
 \mathcal{A} \mid \{ S_i, S_j \}
          IG_i = next(LG_i); IG'_i = new Interval(start = -1)
          IG_{i} = next(LG_{i}); IG'_{i} = new Interval(start = -1)
       elif (IG_i \subset IG_i) / / IG_i disappear during restriction
          IG_{i} = \text{next}(LG_{i}); IG'_{i} = \text{new Interval}(\text{start} = -1)
          if (IG'_{i}[start] = -1) / / IG'_{i}[start] is now known
           IG'_{i}[start] = IG_{i}[start] - shift
       else // (IG_i \subset IG_i) // IG_i disappear during restriction
         ..... // similar to previous case swapping i and j
       shift = shift + |IG_i \cap IG_i|
   elif(IG<sub>i</sub>[start] < IG<sub>i</sub>[start])
     if (IG_i \subset IG_i) / / IG_i disappear during restriction, shift increase
       if (IG/[start] = = -1) // set IG/[start], if not already done, before
 increasing shift
         IG_{i}[start] = IG_{i}[start] - shift
     shift = shift + |IG_i \cap IG_j|
     IG_j = next(LG_j); IG'_i = new Interval(start = -1)
   else // IG_i start after IG_i and is not included in IG_i
     if(IG'_{i}[start] = = -1)
       IG_{i}[start] = IG_{i}[start] - shift
     if (IG_i \cap IG_j \neq \emptyset) / IG_j [start] is now known and shift increase
       IG'_{i}[\text{start}] = IG_{i}[\text{start}] - \text{shift}
       shift = shift + |IG_i \cap IG_j|
     IG'_{i}[end] = IG_{i}[end] - shift
     append IG'_i to LG'_i
     IG_i = next(LG_i); IG'_i = new Interval(start = -1)
```

```
else // (IG_j[start] < IG_i[start])

...... // similar to previous case swapping i and j

if (IG_i \neq NULL) // handle last gaps in LG_i

if (IG_i[start] = = -1) \{IG_i[start] = IG_i[start] - shift\}

IG_i[end] = IG_i[end] - shift

append IG_i to LG_i'

while ((IG_i = next (LG_i) \neq NULL))

append new Interval (start = IG_i[start] - shift; end = IG_i[end] - shift) to LG_i'

if (IG_j \neq NULL)

...../similar to previous block replacing i with j

return LG_i', LG_i'
```

Conclusion

This paper introduces an optimized algorithmic solution to estimate SP-score for affine gap costs in O(nL) and a more versatile solution able to handle more gap cost penalty functions in $O(nL + n^2G_{max})$, with $G_{max} \leq \lceil L/2 \rceil$ being the maximum number of gap intervals per sequence. These optimizations will obviously be part of the next release of MACSE [10], the MSA software we developed to align nucleic sequences with respect to their amino acid translation while allowing them to contain frameshifts and/or stop codons (http://bioweb.supagro.inra.fr/macse/). Moreover, once stated those two algorithms are quite straightforward and can easily be included in the numerous existing MSA software relying on SP-score.

Algorithm. 4. Efficient O(nL) computation of the contribution of gap opening cost (SP_{go}) for an alignment A of n sites and L sequences.

```
<u>Algorithm 4</u>: compute_SP<sub>go</sub>_using_gap_intervals
Input:-the n aligned sequences S_1 \dots S_n of \mathcal{A}
         - the costs of a gap opening within a sequence (gap<sub>0</sub>) or at its extremi-
ties (qap_{0 ext})
Output: SP<sub>go</sub>: the part of the SP-score of \mathcal A due to gap opening costs
nbOpenGap = new Array of L integers initialized to 0;
gapClosing = new Array of L empty lists of Intervals;
for S_i in S_1 \dots S_n
  //construct LG_i in O(L) and update nbOpenGap and gapClosing arrays
  LG_i = \text{compute gap intervals}(S_i)
  foreach IG_i of LG_i
    for k in IG<sub>i</sub>[start] ... IG<sub>i</sub>[end]
      nbOpenGap[k] ++
  append IG<sub>i</sub> to gapClosing[IG<sub>i</sub>[end]]
SP<sub>go</sub> = 0; // part of the SP score related to gap opening costs
for i in 1 ... L
  foreach IG<sub>i</sub> in gapClosing[i]
    if (i == LOR IG_i[start] == 1)
      SP<sub>qo</sub> = SP<sub>qo</sub>+(n-nbOpenGap[IG<sub>i</sub>[start]])<sup>*</sup> gap<sub>0 ext</sub>
    else
      SP_{qo} = SP_{qo} + (n-nbOpenGap[IG_i[start]])^* gap_0
  foreach IG<sub>i</sub> in gapClosing[i]
    for k in IG<sub>i</sub>[start] ... IG<sub>i</sub>[end]
      nbOpenGap[k] = nbOpenGap[k] -1;
return SP<sub>go</sub>
```

Acknowledgments

This work has been supported by the French research agency: Agence Nationale de la Recherche (ANR-10- BINF-01-02 "Ancestrome").

Author Contributions

Wrote the paper: VR.

Conceived and wrote the algorithms: VR.

References

- Przybylski D, Rost B. Alignments grow, secondary structure prediction improves. Proteins. 2002; 46 (2):197–205. PMID: <u>11807948</u>.
- Loytynoja A, Goldman N. Phylogeny-aware gap placement prevents errors in sequence alignment and evolutionary analysis. Science. 2008; 320(5883):1632–5. doi: <u>10.1126/science.1158395</u> PMID: <u>18566285</u>
- Meredith RW, Gatesy J, Murphy WJ, Ryder OA, Springer MS. Molecular decay of the tooth gene Enamelin (ENAM) mirrors the loss of enamel in the fossil record of placental mammals. PLoS Genet. 2009; 5 (9):e1000634. doi: 10.1371/journal.pgen.1000634 PMID: 19730686
- Wong KM, Suchard MA, Huelsenbeck JP. Alignment uncertainty and genomic analysis. Science. 2008; 319(5862):473–6. doi: <u>10.1126/science.1151532</u> PMID: <u>18218900</u>
- 5. Blackburne BP, Whelan S. Class of multiple sequence alignment algorithm affects genomic analysis. Mol Biol Evol. 2013; 30(3):642–53. doi: <u>10.1093/molbev/mss256</u> PMID: <u>23144040</u>
- Altschul SF. Gap costs for multiple sequence alignment. J Theor Biol. 1989; 138(3):297–309. PMID: 2593679
- Wang X-D, Liu J-X, Xu Y, Zhang J. A Survey of Multiple Sequence Alignment Techniques. In: Huang D-S, Bevilacqua V, Premaratne P, editors. Intelligent Computing Theories and Methodologies: 11th International Conference, ICIC 2015, Fuzhou, China, August 20–23, 2015, Proceedings, Part I. Cham: Springer International Publishing; 2015. p. 529–38.
- Edgar RC. MUSCLE: a multiple sequence alignment method with reduced time and space complexity. BMC Bioinformatics. 2004; 5:113.Kececioglu JD, Zhang W. Aligning Alignments. CPM '98: Proceedings of the 9th Annual Symposium on Combinatorial Pattern Matching. Lecture Notes In Computer Science. 1448: Springer-Verlag; 1998. p. 189–208.
- Kececioglu JD, Zhang W. Aligning Alignments. CPM '98: Proceedings of the 9th Annual Symposium on Combinatorial Pattern Matching. Lecture Notes In Computer Science. 1448: Springer-Verlag; 1998. p. 189–208.
- Ranwez V, Harispe S, Delsuc F, Douzery EJ. MACSE: Multiple Alignment of Coding SEquences accounting for frameshifts and stop codons. PLoS One. 2011; 6(9):e22594. doi: <u>10.1371/journal.pone.</u> 0022594 PMID: <u>21949676</u>
- Henikoff S, Henikoff JG. Amino acid substitution matrices from protein blocks. Proc Natl Acad Sci U S A. 1992; 89(22):10915–9. PMID: <u>1438297</u>