

ProbCons: Probabilistic consistency-based multiple sequence alignment

Álvaro Huertas García
Diego Mañanes Cayero
Alejandro Martín Muñoz
Sara Dorado Alfaro

January 16, 2020

Table of contents

- 1 Introduction
- 2 Consistency-based methods
- 3 ProbCons: The algorithm
- 4 Experiments
- 5 Example
- 6 Drawbacks and Conclusions
- 7 References

Table of contents

- 1 Introduction
- 2 Consistency-based methods
- 3 ProbCons: The algorithm
- 4 Experiments
- 5 Example
- 6 Drawbacks and Conclusions
- 7 References

- Multiple sequence alignment (MSA) → way of identifying and visualizing patterns of sequence conservation. It facilitates evolutionary and phylogenetic studies. There are many approaches to multiple sequence alignment:
 - ① Exact methods.
 - ② Progressive alignment (e.g., ClustalW).
 - ③ Iterative approaches (e.g., PRALINE, IterAlign, MUSCLE).
 - ④ Consistency-based methods (e.g., MAFFT, ProbCons).
 - ⑤ Structure-based methods: include information about one or more known 3D protein structures.

Introduction: method's approaches

- Dynamic programming → too inefficient for more than a few sequences. Instead, heuristic strategies: tree-based progressive alignment, sequences are assembled via several pairwise alignment steps. Errors at early stages propagate and may increase the likelihood of misalignment (alleviated by post-processing steps).
- Consistency-based techniques → use evidence from intermediate sequences to guide the pairwise alignment (adjusting the score for a residue pairing according to support from the position of a third sequence that aligns to the others). That is, multiple sequence information is used, as it is being generated.
- COFFEE (another consistency-based) → a library is computed by merging consistent CLUSTALW global and LALIGN local pairwise alignments to form three-way alignments, which are assigned weights. The score for the pairwise alignment is the sum of the weights of all alignments in the library containing that aligned residue pair.

Table of contents

- 1 Introduction
- 2 Consistency-based methods
- 3 ProbCons: The algorithm
- 4 Experiments
- 5 Example
- 6 Drawbacks and Conclusions
- 7 References

Consistency-based methods [SK09]

- Based on: “prevention is the best medicine”
- Combines iterative and progressive approaches with probabilistic models:
 - ① Uses **Hidden Markov Models** to calculate matrices for matching residues in pairwise alignments.
 - ② Uses information about multiple sequence alignment as it is being generated to guide the pairwise alignments.
 - ③ Multiple alignment via tree-based **progressive alignment**
 - ④ Errors at early stages in the alignment are alleviated by **post-processing steps** such as iterative refinement. See [WBH05].

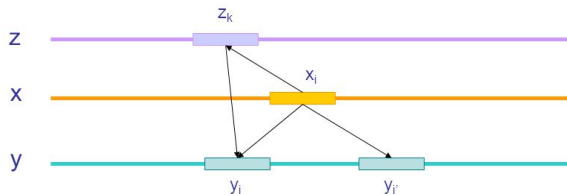
Consistency-based methods

Imaging this biological scenario [Pev15]

Sequence $x \rightarrow x_i$

Sequence $y \rightarrow y_i$

Sequence $z \rightarrow z_k$



- If x_i aligns with z_k and z_k aligns with y_i , then x_i should align with y_i
- Consistency-based techniques **score pairwise alignments** in the context of **information about multiple sequences**

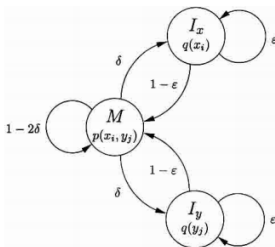
Table of contents

- 1 Introduction
- 2 Consistency-based methods
- 3 ProbCons: The algorithm**
- 4 Experiments
- 5 Example
- 6 Drawbacks and Conclusions
- 7 References

- **ProbCons [DMBB05]** is a pair-hidden Markov model-based progressive alignment algorithm that differs from most typical approaches in its use of **maximum expected accuracy** rather than Viterbi alignment, and of the probabilistic consistency transformation to incorporate multiple sequence conservation information during pairwise alignment.
- **Hidden Markov Models (HMMs)** in sequence analysis are based on a strong probabilistic model that **includes a representation of INDELs** (insertions and deletions, i.e. gaps) [DEKM98].
- The HMM describing families of related sequences are called **profile HMMs**.

Algorithm overview

- In profile HMMs [DEKM98] residues in each position of the alignment can be in one of three possible states:
 - ① **Match:** represent conserved position
 - ② **Insert:** represent small stretches of nonspecific sequence
 - ③ **Delete:** correspond to gaps and represent the absence of a conserved residue
- Each state has associated:
 - ① **Emission probability:** correspond to the probability of observing each amino acid at that particular position of the alignment
 - ② **Transition probability:** describes the frequency of observing a match, insertion or deletion in column $i+1$ given the state column i .



- Emission probabilities, which correspond to traditional substitution scores, are based on the BLOSUM62 matrix.
- Transition probabilities, which correspond to gap penalties, are trained with unsupervised Expectation-Maximization (EM)
 - π_{insert} : initial insertion probability parameter
 - δ : insertion start probability parameter
 - ϵ : insertion extension probability parameter
- The resulting parameters ($\delta = 0.019931$, $\epsilon = 0.79433$, $\pi_{insert} = 0.19598$) are used as default by the program.

ProbCons [DMBB05]

- Given m sequences $\rightarrow S = \{s^{(1)}, \dots, s^{(m)}\}$.
- Maximum expected accuracy.
- Probabilistic consistency \rightarrow MSA conservation information in the pairwise alignment.

- 1 Step 1: Computation of posterior probability matrices.
- 2 Step 2: Computation of expected accuracies.
- 3 Step 3: Probabilistic consistency transformation.
- 4 Step 4: Computation of the guide tree.
- 5 Step 5: Progressive alignment.
- 6 Step 6: Iterative refinement (post-processing OPTIONAL step).

Step 1: Computation of posterior probability matrices

- For $x, y \in S$, compute the matrix

$$P_{xy}(i, j) = \mathbf{P}(x_i \sim y_j \in a^* | x, y) ,$$

where $1 \leq i \leq |x|$ and $1 \leq j \leq |y|$.

- Each position $P_{xy}(i, j)$ is the **posterior** probability that letters x_i and y_j are paired in an alignment a^* .
 - Computing posterior probabilities in pair-HMMs [DEKM98].
- Time complexity $O(m^2 L^2)$.
 - m is the number of sequences.
 - L is the length of each sequence.

Step 2: Computation of expected accuracies

- The expected accuracy is defined as

$$\mathbf{E}_{a^*}(\text{acc}(a, a^*)|x, y) = \frac{1}{\min\{|x|, |y|\}} \sum_{x_i \sim y_j \in a} P_{xy}(i, j) ,$$

where a is the alignment that maximizes the expected accuracy by dynamic programming.

- Set

$$E(x, y) = \mathbf{E}_{a^*}(\text{acc}(a, a^*)|x, y) . \quad (1)$$

Step 3: Probabilistic consistency transformation

- Reestimate quality scores $\mathbf{P}_{xy} \rightarrow$ probabilistic consistency transformation.
- Incorporate similarity of x and y to other sequences in S :

$$\mathbf{P}'(x_i \sim y_j \in a^* | x, y) = \frac{1}{|S|} \sum_{z \in S} \sum_{z_k \in Z} F(x_i, y_j, z_k),$$

where $F(x_i, y_j, z_k) = \mathbf{P}(x_i \sim z_k \in a^* | x, z) \times \mathbf{P}(z_k \sim y_j \in a^* | z, y)$.

- In matrix form:

$$\mathbf{P}'_{xy} = \frac{1}{|S|} \sum_{z \in S} \mathbf{P}_{xz} \mathbf{P}_{zy}.$$

- **Optimization:** use sparse matrices ignoring entries $\leq \omega$ (threshold).
- This step can be iterated until convergence.

Steps 4, 5 and 6

- Hierarchical clustering.
 - Similarity measure $E(x, y)$ as defined in Equation (1).
 - WPGMA method.
- Align sequence groups hierarchically.
 - Sum-of-pairs.
 - Gap penalties $\rightarrow 0$.
- Progressive alignment [WBH05].
 - Randomly partition alignment into two groups of sequences.
 - Realign.
 - This step can be iterated.

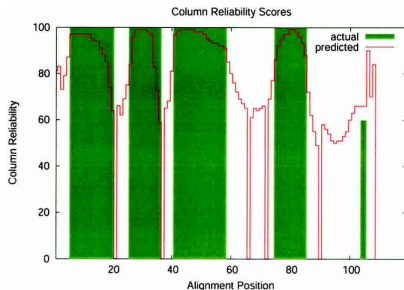
Table of contents

- 1 Introduction
- 2 Consistency-based methods
- 3 ProbCons: The algorithm
- 4 Experiments**
- 5 Example
- 6 Drawbacks and Conclusions
- 7 References

Some experiments with BALiBASE dataset

- The BALiBASE dataset:
 - 141 reference protein alignments.
 - Hand constructed alignments from the literature.
 - 5 subsets with alignments of different characteristics.
 - Test alignments are scored respect **core blocks** → reliable alignments.
- **No universally accepted accuracy measure for protein alignments.**
 - Sum-of-pairs score (SP).
 - Column score (CS).

Column reliability for BALiBASE



```

1csy SHEKMPWFHGKISRREEQIVLGSKTNGKFLIRARD--NNGSYALCLLNEGKVLHYRID
1gri EMKPHPWFFGKI PRAKAEDML-SKQRHGDGAFIRESE-SAPGDFLSVKFGNDVQHFKVL
1aya ---MRRWFHFNITGVAEHLL-LTRGVGGSFLARPSK-SNP GDFLSVRENGAVTHIKIQ
2pna -LQDAEFTWGDISRREEVNEKL--RDTADGTFVLVDASTPMHGDDYTLTLREGGNSKLIKIF
1bfi HHDEKTFNVGSSNRNKAENLL--RGKR DGTFLVRESS--KQGCYACSVVVDGEVHCVIN
      1      11      21      31      41      51
    
```

```

1csy KDKTKGLSIEEG-KKPTLMQLVEHYSYKA-----DGLLRVLTVFCQK
1gri RDGAGK-YFLAV-VKFNLSLNLVDYHRSTS-VSRHQQIFLRDIEQVTFQK
1aya NTGDIY-DLYGG-EKFATLAEVLQYIMEHHGQKKEKNGDVIELKYP-LN
2pna HRDGKY-GFSDP-LTFNSVVELINHYRNES-LAQYNPKLDVKLLYP-VS
1bfi KTATGY-GFAEPTNLYSSLKELVLYHQHSTS-LVQHNDLSLNTLAYFVYA
      61      71      81      91     101
    
```

Image from [DMBB05].

At each position:

- Red line → predicted proportion of correct pairwise matches.
- Green Blocks → actual proportion of correct pairwise matches.

Comparison with other methods

ProbCons multiple alignment tool

Table 1. Performance of aligners on the BALiBASE benchmark alignments database

Aligner	Ref 1 (82)		Ref 2 (23)		Ref 3 (12)		Ref 4 (12)		Ref 5 (12)		Overall (141)		Time (mm:ss)
	SP	CS	SP	CS	SP	CS	SP	CS	SP	CS	SP	CS	
Align-m	76.6	n/a	88.4	n/a	68.4	n/a	91.1	n/a	91.7	n/a	80.4	n/a	19:25
DIALIGN	81.1	70.9	89.3	35.9	68.4	34.4	89.7	76.2	94.0	84.3	83.2	63.7	2:53
CLUSTALW	86.1	77.3	93.2	56.8	75.3	46.0	83.4	52.2	85.9	63.8	86.1	68.0	1:07
MAFFT	86.7	78.1	92.4	50.2	78.8	50.4	91.6	72.7	96.3	85.9	88.2	71.4	1:18
T-Coffee	86.6	77.4	93.4	56.1	78.5	48.7	91.8	73.0	95.8	90.3	88.3	72.2	21:31
MUSCLE	88.7	80.8	93.5	56.3	82.5	56.4	87.6	60.9	96.8	90.2	89.6	73.9	1:05
ProbCons	90.1	82.6	94.4	61.3	84.1	61.3	90.1	72.3	97.9	91.9	91.0	77.2	5:32
ProbCons-ext	90.0	82.5	94.2	59.1	84.3	61.1	93.8	81.0	98.1	92.2	91.2	77.6	8:02

Columns show the average sum-of-pairs (SP) and column scores (CS) achieved by each aligner for each of the five BALiBASE references. All scores have been multiplied by 100. The number of sequences in each reference is given in parentheses. Overall numbers for the entire database are reported in addition to the total running time of each aligner for all 141 alignments. The best results in each column are shown in bold.

Figure: Image from [DMBB05].

Table of contents

- 1 Introduction
- 2 Consistency-based methods
- 3 ProbCons: The algorithm
- 4 Experiments
- 5 Example**
- 6 Drawbacks and Conclusions
- 7 References

Examples: Comparison between methods

- MSA of distantly related globins (human beta globin, human myoglobin, human neuroglobin, soybean leghemoglobin, rice hemoglobin) using four different programs. Symbols: * complete conservation, : conservative substitutions, . less conservative substitutions. Programs differ in:
 - Align corresponding regions of alpha helical secondary structure (red lettering).
 - Align conserved histidines (open and black arrowhead). They are important in coordinating protein binding to the heme group → they should be aligned by all the programs. The open arrowhead histidine shows a complete conservation. The conservation of the black is only achieved by ProbCons and T-Coffee.
 - Create and place gaps (boxed regions).

(a) Praline multiple sequence alignment

```

beta globin      .....MVLHTPEEKSAVTALGVK.V..INVDGGGAELGRLLVVPWTQRFPES.PG
myoglobin       .....MGLSDGQWGLNVLWVNGKVEADI.PHQGGQVILRLFKLPEPTLEKDFP.KK
neuroglobin     .....MERPEPLRIGQNRVAISR.PLEHTGLVFLARLFGEPDLLPLFLQVCR
soybean         .....MVAFTPKQDQALVSSSFAFKANI.PQSVVFVYTSLLEKAPAKQDLS...FL
rice            MALVEDINNAVVASPSLEQELVLKSAI.LKDSVIANLRFPFLKFLVPEVASPSAQS
Consistency     0000000000142654382579345734636364343624538664333534450063

beta globin      DLKSTPDVAVNGNPKDKKVGKGLVGFADSG.AHLNDLNGTFLATSLSL..NCKDLK...VDP
myoglobin       HLGSEDMKMAKEDKDKGHATVATLALOGI.LKKGKHGEARICPAQS..KATHKK...IPV
neuroglobin     QFSSEPCDLSSEFLDHRKRVMLVIDAATPVNEDLSLSEYLSLGRKRGK...VKL.A
soybean         A.MGVDP..TNPKLTHGAELFALVRDSAGQLL.KASGVTFVDA...LGSVAGKAVAT
rice            R.NSDVPELKKKLTTHAMVFWCCEAAV.LKAL.RKAGVTVTRDQKSLGATLTKYGVGD
Consistency     3163454224776653436863552424644543536343335420033544000922

beta globin      ENPLRLGNLVQVYLAHIF..GKPEFTPVQAYQKQVAGVANAIAHKHYH...
myoglobin       KYLFIEISCTIYVCLLQGL..PGDFGADGAGKMGVAFLEPDMASMLKELPQSG
neuroglobin     SPSFTVSESLIYNLYKEGL..GPAPTATRAWQLYGVAGVGVNSRGMD..GE
soybean         PPTVVVVEALKLTIAAY..GDHSELSRGEAVYDLVDAEALAKKA...
rice            AHFVEVVKFALLOLTIEKVPDMSKSPAMKSAWEAYDHLVLAIGKMPKAE...
Consistency     4374484449482854523053365544545546562644675342001000

```

(c) PROBCONS

```

beta globin      M-----VHLTPEEKSAVTLAKGVNVAD-----EVGGELAGRLVVVYPWTQPFES-PG
myoglobin       M-----GLDGQWGLVLMVVGKVRVADIPGHQGVLRITLFGKHPTLEKDFK-FK
neuroglobin     M-----ERPEPELRIGQSRVASSRSPLEHQHTVLFARLFALEPDLFLPQVNCR
soybean         M-----VAFTEKQDALVSSPEAFKNIIPKVVVPTVTSITPKSAAKDLI-FK-IA
rice            MALVEDNNNAVSESEGEALVLSKWAALKKDSANIALRPFLLKTFEVPASASQMPF-LR
                *      *      *      *      *      *      *      *      *      *
                *      *      *      *      *      *      *      *      *      *
beta globin     DLSTPDVAMGNPKVKAAGKGVKGLVAFSDGVAHLHD---NLK---GTFATLSRLAKDLHVDP
myoglobin       HKSKDEKMAKEDLKKGHATVITLALGQI---LKKKGHHH---ARIKPLAKSKTHKVP-FG
neuroglobin     QFSSDEPCLDSSPELIRKIVMLVDATAVTVNEDLSLE---EYLASLGRKRRVA-QVKL
soybean         NGVNDP---TNPKLTGHAELPALVRDSSQQLKAGSTV---ADAALGSVVAQK-A-VD
rice            NSDVP---LEKNPKLTKHMSVPMTCIAAQLRKAGVTTVRDTLLKGLGATLKY-GTD
                *      *      *      *      *      *      *      *      *      *
                *      *      *      *      *      *      *      *      *      *
beta globin     ENFRLLGNGLVLIQVHFG-KGEFTPPVQAQVGVVAGVANAIAHK-----YH
myoglobin       KYLEFISECIIVLQSKH-PDGFDAGAGAMMKALEFKRIMASNYKELGPGF
neuroglobin     SSSFSTGEALLMYLEKLCI-GPAFTPATRMASSVLYGAVGVAMSRG---N-DGE
soybean         PQFVVVKEALDKTKIAEVA-GDRVSDLSRAMEVVEDELAIAAK---KA
rice            AHFEVVKFALLDTKEEVPADMWSPAMKSMASVEYDHLVAAIKQE---MKPAE

```

(b) MUSCLE (3.6) multiple sequence alignment

```

beta globin      -----MVHLTPEEKSAVTALMGKVNVD-----EVGGEALGRLLVVYPWTQRFES-FG
myoglobin       -----MGLSDGEQLVILNKGKVEADIPGHQEVYLILRFGKHPETLEKFK-FK
neuroglobin     -----MERPEPEILRQSRWRAVEDSEHGTEVTLFARLFALFEPOLLPLFYQICR
soybean         -----MVAFTEKQDALVSSFEAFKANIPQYSVVYFYSILEKAPAKDLFGFFLA
rice            MALVEDNNNAVAVGSEEQEALVLKLSWALKDKDANSALRFFTKLFIKFEVPASGAFSG-LR

                                     *      *      *
beta globin      DLSTPDVAMGNPKVAGSGKGVLGAF-----SDGLAHLNLDKATGLTSELKSLGDKLIL-VDPE
myoglobin       HLKSKEDMSAESEDLKKGIAVLTAL-----GGILKKKHGHEARLSEKQASNAIKH-IPVK
neuroglobin     NSSPEEDLSPPSLSEHILKVMVL-----DAATNVNEDLSLSEVLEGKRRARAVKVL
soybean         NGVDSP-----TNPLKLTHGAEKLFALVRDASGKASGTSTVYD-----AALGSVGAQNGKTS
rice            NSDVP-----LEKSNPKLKTTRANSVYVMTCEAAAGKAGKGVTVRDITTLKRIGATHLFGVGDGA

                                     *      *      *
beta globin      NFRLRGVLCVLIHLHGDEE-FTPPVQAAYRKVVAAGVANALAHYH-
myoglobin       YLFSEITSEI QIVQLSGKPGD-FGADQAAGNMRALKFLKCDMSASYRGKFGQ
neuroglobin     FPGTSGVSSLLITMLEKCLGPA-FTPATRAANSQLYGAVANSSYRGKGDGE
soybean         FGVVVKAEALIT KAAVQD-WSEDLRSNRAVEYD ELAAAIKKA
rice            HF EVVKPALLOT LK EEPVALMWSAPMSK SANSEAYDHLVAIK QEMKPAE-

```

(d) CLUSTAL FORMAT for T-COFFEE Version 5.13

```

beta globin      -----MVHLTPTEEEKSAVLNMGKGVND--EVGVEGLGRLLVYWPWTQTEFDF--SPG
myoglobin       -----MGLSDGEQWLQVLMGKGVKVD+PQHGGVRLILFKPHGPIETLEKDF-KFK
neuroglobin     -----MERPEDELRQSRVAVRSSPFLVYFARLFALEFADLFLDLPQYKNC
soybean         -----MVAFETKQDALVSSSEAFKANI+POYSVVFYTSILEKAPAKDLFS--FLA
rice            MALVEDNNAVAVSFSEEQEALVLKSNAILKDDSANIALRFFPKIFVFPASQMFPS--FLR

              *      *      *      *      *      *      *      *      *      *
              |      |      |      |      |      |      |      |      |      |
              v      v      v      v      v      v      v      v      v      v
beta globin     DLSTDPDAVMGNPKYKAHKGVQLGAPSDGLAHLNDLNL---KGTFF---ATLSLEKADLKHVP
myoglobin       HLKSEDEMCASEDILKFGDAGTVTLA---GGILKKKGHEARE---TKPLAQSHAKHKIPV
neuroglobin     QFSSPEDCLSSPEKPHDRIKVLMLVIDAA+TNVEDNL---SSLEYLALSLGRKH--RAQVGLT
soybean         NGVDNP---TNPKITGAHAKFALVRDSSAGGLKASGTVYND---AALSGVYKAGKATVTP
rice            NSDVPV---LEKNPKIKTHAMSVFVMTCEAAAGLRKAKGVTYVTDITKKLGATHLKTGVGDLA

              *      *      *      *      *      *      *      *      *      *
              |      |      |      |      |      |      |      |      |      |
              v      v      v      v      v      v      v      v      v      v
beta globin     ENFRLLGNVLVCLVLAHFG--GKEFTTPVQAAYQKVVAQVANALAHKYH---
myoglobin       KSLFESBICIQLVQSKH--PGDFGADAGQAMNKALEFLFKDMKASHIRGLGPGQ
neuroglobin     SYFSTVSGISYLMYQSKL(GPAFTTPATRAANSLQIYAVVQAMSHRWGQD---E
soybean         Q-FVVVEKALLATIKAAV--GKWSDELSRANVEYDELAJAAIKQA---
rice            H--FEVFKFALLDTIKAEVVPDMMSPAMKSAWSEATHLVLAALQK---MKPAP

```


Table of contents

- 1 Introduction
- 2 Consistency-based methods
- 3 ProbCons: The algorithm
- 4 Experiments
- 5 Example
- 6 Drawbacks and Conclusions**
- 7 References

- **Computational weight:** The computation step of calculation of posterior probabilities takes time $O(m^2L^2)$, where m is the number of sequences and L is the length of each sequence.
- **M-Coffe (Meta-Coffe):** combines the output of 15 different sequence alignment methods(ProbCons included). M-Coffe employs a consistency-based approach to estimate a more accurate consensus alignment.
- **Structural methods:** adding structural information, even further accuracy is achieved.

Conclusions

The ProbCons algorithm uses an extremely simple model of sequence similarity (a three-state pair-HMM):

- 1 Makes no attempt to incorporate biological knowledge (i.e position specific gap scoring or rigorous evolutionary tree construction).
- 2 Use amino acid alphabet and BLOSUM emission probability matrices as protein-specific alignment information.
- 3 Can be used to DNA alignment by changing the alphabet and the BLOSUM matrices with values for nucleotides.
- 4 The parameter used in the model are transparent (π_{insert} , δ , ϵ).
- 5 The training program is done automatically on unaligned sequences using Expectation-Maximization.
- 6 High accuracy: probabilistic consistency transformation and objective function.

Table of contents

- 1 Introduction
- 2 Consistency-based methods
- 3 ProbCons: The algorithm
- 4 Experiments
- 5 Example
- 6 Drawbacks and Conclusions
- 7 References**

References I



Richard Durbin, Sean R Eddy, Anders Krogh, and Graeme Mitchison, *Biological sequence analysis: probabilistic models of proteins and nucleic acids*, Cambridge university press, 1998.



Chuong B Do, Mahathi SP Mahabhashyam, Michael Brudno, and Serafim Batzoglou, *Probcons: Probabilistic consistency-based multiple sequence alignment*, Genome research **15** (2005), no. 2, 330–340.



Jonathan Pevsner, *Bioinformatics and functional genomics*, John Wiley & Sons, 2015.



Fahad Saeed and Ashfaq Khokhar, *An overview of multiple sequence alignment systems*, arXiv preprint arXiv:0901.2747 (2009).



Iain M Wallace, Gordon Blackshields, and Desmond G Higgins, *Multiple sequence alignments*, Current opinion in structural biology **15** (2005), no. 3, 261–266.

ProbCons: Probabilistic consistency-based multiple sequence alignment

Álvaro Huertas García
Diego Mañanes Cayero
Alejandro Martín Muñoz
Sara Dorado Alfaro

January 16, 2020