

ppALIGN: Posterior distribution for score-based alignments

Stefan Wolfsheimer
Gregory Nuel

Mathématiques Appliqués à Paris 5
Université Paris Descartes



Outline

Score-based alignment

Probabilistic alignment: ppALIGN

ppALIGN in action

What is an alignment ?

Example (Two DNA sequences)

$$\begin{aligned}a_1^\ell &= \text{a c g t a g c a t g a c a} \\ b_1^m &= \text{a c c g t a c a a g c a}\end{aligned}$$

c: common ancestral sequence

c	a	c	c	g	t	t	a		c	a	a	g	a	c	a
a_1^ℓ	a	c		g	t		a	g	c	a	t	g	a	c	a
b_1^m	a	c	c	g	t		a		c	a	a	g		c	a

here is the **alignment** $\mathcal{A} = (\tilde{a}_1^t, \tilde{b}_1^t)$ we get:

\tilde{a}_1^t	a	c	-	g	t		a	g	c	a	t	g	a	c	a
\tilde{b}_1^t	a	c	c	g	t		a	-	c	a	a	g	-	c	a

Note that $\max(\ell, m) \leq t \leq \ell + m$.

Score-based alignments

Definition (Score of an Alignment)

- **scoring function:** $\sigma : \Sigma \cup \{-\} \times \Sigma \cup \{-\} \rightarrow \mathbb{R}$
- **score of an alignment:** $s(\mathcal{A}) = s(\tilde{\mathbf{a}}_1^t, \tilde{\mathbf{b}}_1^t) = \sum_{k=1}^t \sigma(\tilde{\mathbf{a}}_k, \tilde{\mathbf{b}}_k)$
- optimal score $s_0 = \max_{\mathcal{A}} s(\mathcal{A}; \mathbf{a}_1^\ell, \mathbf{b}_1^m)$ and alignment $\mathcal{A}_0 = \operatorname{argmax}_{\mathcal{A}} s(\mathcal{A}; \mathbf{a}_1^\ell, \mathbf{b}_1^m)$
- **Needleman-Wunsch algorithm (1970)**
- Problem in score-based alignment: (nearly) optimal alignments **not unique**.
- Reliable regions: common to all high scoring alignments
- Questionable regions: close to gaps and low complexity regions

Posterior probabilities

Example (protein alignment)

SALLASGGTSSHRWSRT	score = 31
SALLMARKSHRVLWSRT	
<hr/>	
SALLASGGTSSHR--WSRT	score = 31
SALLMA--RKSHRVLWSRT	
<hr/>	
SALLASGGTSSHR--WSRT	score = 28
SALL--MARKSHRVLWSRT	



SALLASGGTSSHRWSRT
 |||||+++++|||
 SALLMARKSHRVLWSRT

score = 31



SALLASGGTSSHR--WSRT
 |||||++ ++||| |||
 SALLMA--RKSHRVLWSRT

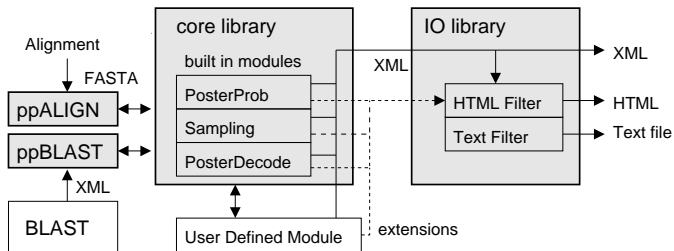
score = 31



SALLASGGTSSHR--WSRT
 ||| ++||| |||
 SALL--MARKSHRVLWSRT

score = 28

ppALIGN



Features

- Alternative decoding algorithms
- Input filter for BLAST
- Structured output (XML)
- Standalone programs and C++ library (open source)
- Webinterface <http://www.math-info.univ-paris5.fr/ppblast/>

SW, AK Hartmann, G Nuel, preprint



Protein alignment

SCORE PARAMETERS

Score Matrix: ?

Gap open: ?

Gap extension: ?

ppALIGN PARAMETERS

Model: ? Pair Hidden Markov Model

? Finite temperature alignment

ALGORITHM

Algorithm: ? Global alignment

? Local alignment

? Boundaries of local alignments

Alignment: ? Find optimal alignment

? User defined alignment

ALTERNATIVE ALIGNMENTS

Decoding method: ? Maximum accuracy alignment

? Sampling from the posterior distribution

Number of sampels:

USER DEFINED ALIGNMENT

Compute the posterior probabilities of the following global alignment.
The number of columns in both sequences must agree.

Alignment (Fasta Format):

```
>62810|unnamed protein product
LQVSTPOPVNAGS--EDESGKG-----NLGFIHAFVASISVIIIVSELGDKTFFIAAIM
AMRYNRLVVLGAMLALGVMTCLSVLFQYATTIIPRIYTYVSTALFAIFGIRMLREGLK
>gi|170036348|ref|XP_001846026.1|conserved hypothetical protein [Culex
quinquefasciatus]
VTELSNPNVESGSPGEEKSSAGGGLSSDVGFMHAFIASFVIIIVSELGDKTFFIAAIM
AMRHPRLTVFAGAIAALALMTVLSAVFGMAATIIIPRVYTYI STALFALFGLKMLKEGY
```

ppALIGN

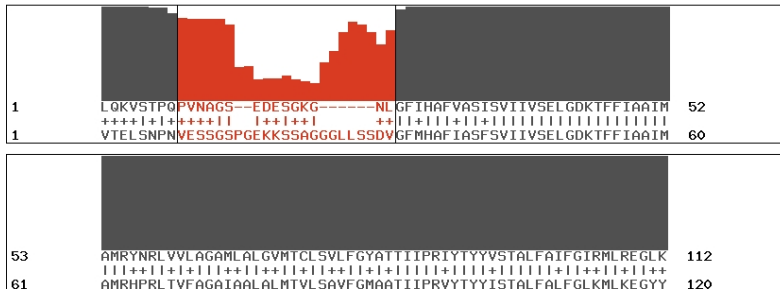
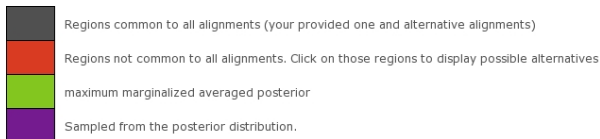
ppALIGN Result

Parameters

[\[Show\]](#)

Alignment

[\[Show Details\]](#)



ppALIGN Result

Parameters

[\[Show\]](#)

Alignment

[\[Show Details\]](#)

