

Stochastic Pairwise Alignment

Ulrike Mückstein

Institute for Theoretical Chemistry and Structural

University Vienna

<http://www.tbi.univie.ac.at/~ulim/>

Bled, 2003

Motivation

- ★ Alignment algorithms normally return a single optimal alignment. However, an optimal alignment does not need to be the only optimal alignment for the optimization problem. Additionally, suboptimal alignments of equal quality might exist.
- ★ The optimal alignment of distantly related sequences may be highly sensitive and susceptible to small perturbations of scoring parameters.
- ★ The reliability of an alignment depends strongly on the degree of sequence similarity.

Approaches dealing with locally variable alignment

Variations of dynamic programming algorithms to construct
alignments

Vingron & Argos, 1990;

Saqi & Sternberg, 1991;

Calculation of the *partition function* and of match probab

Miyazawa, 1994;

Generation of ensembles of suboptimal alignments by s
tracking

this work;

Sequence Alignments

score $s(a, b)$ of a match between two residues

$$s(a, b) = k \log \frac{f_{ab}}{f_a f_b}$$

additive alignment score function $S(\mathcal{A})$

$$S(\mathcal{A}) = S(\mathcal{A}_{1,1}^{i-1,j-1}) + s(a_i, b_j) + S(\mathcal{A}_{i+1,j}^{m,n})$$

Probabilistic Interpretation of Sequence A

probability of a particular alignment \mathcal{A}

$$\text{Prob}(\mathcal{A}) = \frac{1}{Z} \exp\left(\frac{S(\mathcal{A})}{k}\right)$$

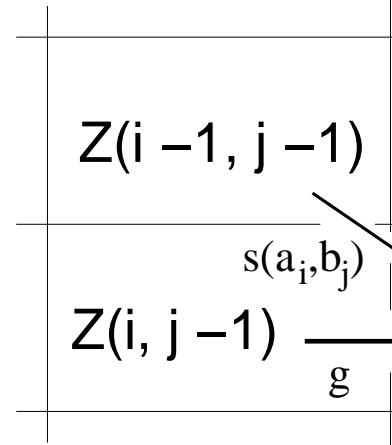
partition function

$$Z = \sum_{\mathcal{A}} e^{\frac{S(\mathcal{A})}{k}}$$

see e.g. Miyazawa, 1994, Yu & Hwa, 2001.

Partition Function

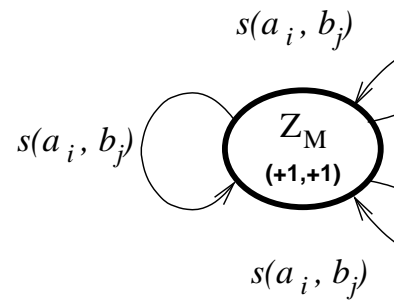
$$\begin{aligned}
 Z_{i,j} &= Z_{i-1,j-1}e^{\beta s(a_i,b_i)} \\
 &\quad + Z_{i,j-1}e^{\beta g} \\
 &\quad + Z_{i-1,j}e^{\beta g}
 \end{aligned}$$



$$\beta = \frac{1}{k}$$

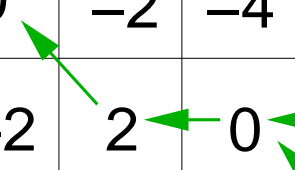
canonical Alignments

A---XXXXB and AXXXX---B
 AYYY----B A----YYVB



Stochastic Backtracking

	-	A	U
-	0	-2	-4
A	-2	2	0
G	-4	0	2



next state	next indices $i \leftarrow$ $j \leftarrow$	transition probability
match	$i - 1$ $j - 1$	$\frac{Z_{i-1,j-1} e^{\beta s(a_i,b_j)}}{Z_{i,j}}$
gap in a	i $j - 1$	$\frac{Z_{i,j-1} e^{\beta g}}{Z_{i,j}}$
gap in b	$i - 1$ j	$\frac{Z_{i-1,j} e^{\beta g}}{Z_{i,j}}$

Match Probabilities

$$P_{ij} = \frac{Z_{ij}^M \widehat{Z}_{ij}^M}{Z} \exp(-\beta s(a_i, b_j))$$

partition function of all alignments

Z_{ij}^M of the sub-sequences $a[1..i]$ and $b[1..j]$ ending with a match of (a_i, b_j)

\widehat{Z}_{ij}^M of the sub-sequences $a[i..m]$ and $b[j..n]$ starting with a match of (a_i, b_j)

see Miyazawa, 1994, Vingron & Argos, 1990.

Comparison with Structure Based Alignment

3D struct. _____
 GALTESQAALVKSSWEEFNANIPKHTRFFILVLEIAPAAKDLFSKLGKTSEVPQNM-PELQAHAG
 -VLSPADKTNVKAAWGKVGAAHAGEYGAEALERMFSLFPPTTKTYF--PHFDLSHGSAQVKGHGKKV

optimal
 GALTESQAALVKSSWEEFNANIPKHTRFFILVLEIAPAAKDLFSFLKGTSEVPQNNPELQAHAG
 -VLSPADKTNVKAAWGKVGAAHAGEYGAEALERMFSLFPPTTKTYFPHF----DLSHGSAQVKGHGK

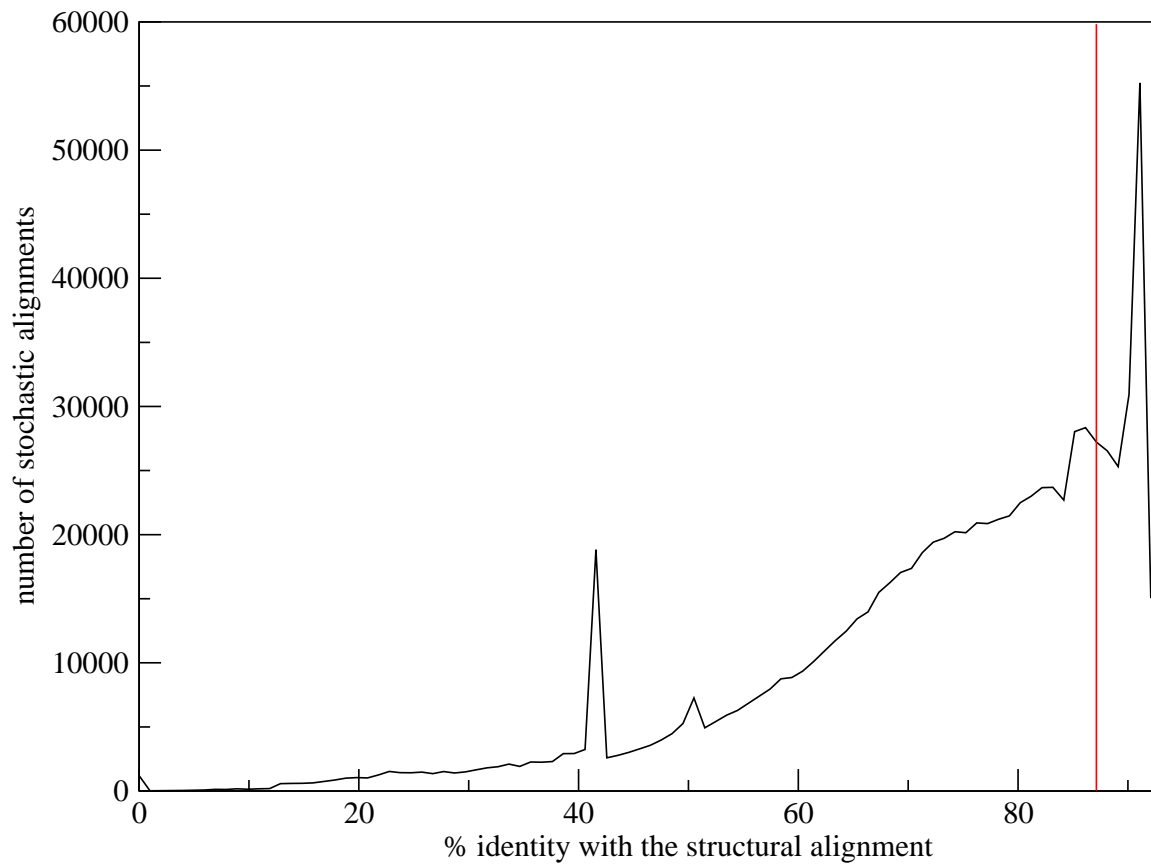
3D struct. _____
 VTGVVVTDATLKNLGSVHV-SKGVADAHFPVVKAEAILKTIKEVVGAKWSEELNSAWTIAYDELA
 ----DDMPNALALSDDLHAHKLRVDPVNFKLLSHCLLVTLAAHLPAEFTPAVHASLDKFLASVST

optimal
 TGVVVTDATLKNLGSVHVSKGVAD--AHFPVVKAEAILKTIKEVVGAKWSEELNSAWTIAYDELA
 MPNALSALSDLHAHKLRVDP-----VNFKLLSHCLLVTLAAHLPAEFTPAVHASLDKFLASVST



alignment of *L.luteus* Leghaemoglobin / human Hemoglobin.

Stochastic Backtracking Finds Suboptimal Alignments that are the Optimal Alignment



Comparison with Structural Alignments by CE

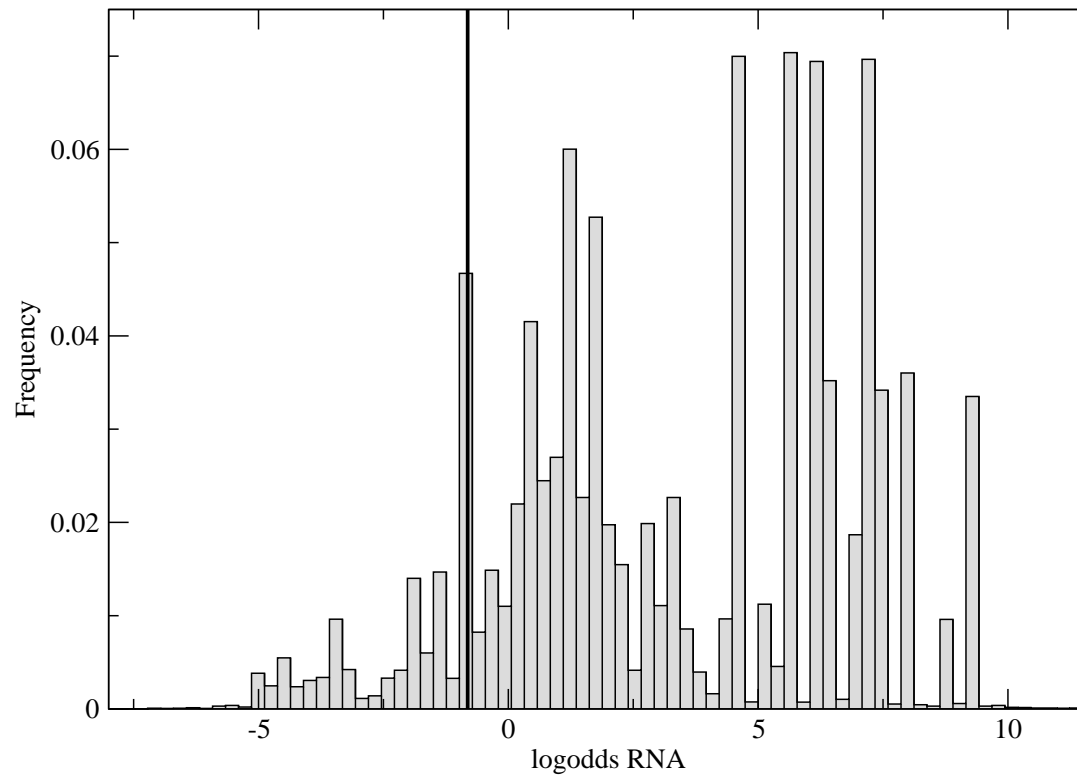
PDB entry	length	% identity		% eq be than
		optimal	mean of stoch. aln.	
1YCC/1CTJ ¹	108/89	51.3	50.3	4.
4HHB/1GDJ ²	141/153	67.2	63.5	38
1SAC/1C4R ³	204/182	44.4	38.4	2.

¹ Cytochrome C, *S.cer.* / Cytochrome C6, *M.braunii*

² Hemoglobin, Human / Leghaemoglobin, *L.luteus*

³ Serum Amyloid P Component, Human / Neurexin (partial), *R.nor.*

Using Stochastic Backtracking with other M



Classification of the TAR hairpin of HIV with `qrna`, Rivas & Ed

Conclusion

- ▶ In this work we present an algorithm that produces correct samples of alignments by stochastic backtracking.
- ▶ The ensemble of stochastic alignments contains correct alignments with significant probabilities even though the optimal alignment differs significantly from the structural alignment.
- ▶ Stochastic pairwise alignments can be used as input data for dynamic programming tools.
- ▶ The software package can be downloaded from the internet at <http://www.tbi.univie.ac.at/~ulim/probA/>.

- ▶ iterative multiple alignment procedures are likely to be trapped in local optima that differ from the true alignment
- ▶ use the match probability matrix of the pairwise alignment of stochastic pairwise alignments to develop a multiple alignment

do not use canonical alignments

A---XXXXB	and	AXXXX---B
AYYY----B		A----YYB
CXXXYYC		CXXXYYC

Thanks to

Peter Stadler, Ivo Hofacker, Peter Schuster.